



# Duality and equivalencies in closed tandem queuing

Zhen Liu, Jacques Resing

## ► To cite this version:

Zhen Liu, Jacques Resing. Duality and equivalencies in closed tandem queuing. [Research Report] RR-2115, INRIA. 1993. inria-00074557

**HAL Id: inria-00074557**

**<https://hal.inria.fr/inria-00074557>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

*Duality and Equivalencies  
in Closed Tandem Queueing Networks*

Zhen LIU  
Jacques RESING

N° 2115  
Septembre 1993

PROGRAMME 1

Architectures parallèles,  
bases de données,  
réseaux et systèmes distribués

*Rapport  
de recherche*

1993

# Duality and Equivalencies in Closed Tandem Queueing Networks

Zhen Liu\*

INRIA Centre Sophia Antipolis  
2004, route des Lucioles  
B.P. 93, 06902 Sophia Antipolis  
France

Jacques Resing†

Department of Mathematics and Computing Science  
Eindhoven University of Technology  
P.O.Box 513  
5600 MB Eindhoven  
The Netherlands

September 15, 1993

---

\*The work of this author was supported in part by the CEC DG XIII under the ESPRIT BRA grant QMIPS.

†The research leading to this paper has been carried out while this author was visiting INRIA, Centre Sophia Antipolis. This visit has been made possible by financial support of the Netherlands Organization for Scientific Research (NWO). Kind hospitality of the MEVAL group at INRIA during this visit is also greatly acknowledged.

# Duality and Equivalencies in Closed Tandem Queueing Networks

## Abstract

Equivalence relations between closed tandem queueing networks are established. Four types of models are under consideration: single-server infinite-capacity buffer queues, infinite-server queues with resequencing, single-server unit-capacity buffer queues with blocking before service, and single-server unit-capacity buffer queues with blocking after service. Using a customer/server duality we show that in a network consisting of single-server infinite-capacity-buffer queues, customer-dependent service times and server-dependent service times yield equivalent performance characteristics. We further show that for closed tandem queueing networks, a system consisting of single-server infinite-capacity-buffer queues (resp. infinite-server queues with resequencing) and a system consisting of single-server unit-capacity-buffer queues with blocking before service (resp. blocking after service) have equivalent performance behaviors. As applications of these equivalence properties, we obtain new results on the analysis of symmetric closed tandem networks, where all the service times are independent and identically distributed random variables. In particular, we obtain a closed-form expression for the throughput of networks with unit-capacity-buffer queues and blocking before service when the service times are exponentially distributed. We also prove the monotonicity of throughput (of queues) with respect to the number of queues and number of customers in these models. This last property in turn implies the existence of nonzero asymptotic throughput when the number of queues and number of customers go to infinity.

**Keywords:** Equivalence, customer/server duality, closed tandem queueing network, infinite server with resequencing, queues with blocking, customer-dependent service, symmetric queueing network.

# Dualité et équivalences dans les réseaux de files d'attente en tandem fermés

## Résumé

Des relations d'équivalence entre des réseaux de files d'attente en tandem fermés sont établies. Quatre types de modèles sont analysés : files d'attente à serveur simple et tampons de capacité infinie, files d'attente à serveur infini avec reséquencement, files d'attente à serveur simple et tampons de capacité unitaire avec blocage avant service, files d'attente à serveur simple et tampons de capacité unitaire avec blocage après service. Utilisant une dualité de client/serveur, nous démontrons que dans un réseau de files d'attente à serveur simple et tampons de capacité infinie, les temps de service dépendant des clients et les temps de service dépendant des serveurs fournissent des caractéristiques de performances équivalentes. Nous démontrons ensuite que pour ces réseaux de files d'attente en tandem fermés, un système de files d'attente à serveur simple et tampons de capacité infinie (resp. files d'attente à serveur infini avec reséquencement) et un système de files d'attente à serveur simple et tampons de capacité unitaire avec blocage avant service (resp. files d'attente à serveur simple et tampons de capacité unitaire avec blocage après service) ont des performances équivalentes. Comme applications de ces propriétés d'équivalence, nous obtenons nouveaux résultats sur l'analyse des réseaux symétriques où les temps de service sont des variables aléatoires i.i.d. En particulier, nous obtenons une expression sous forme close pour le débit des réseaux avec tampons de capacité unitaire et blocage avant service quand les temps de service ont une distribution exponentielle. Nous prouvons aussi la monotonie du débit (des files d'attente) par rapport aux nombre de files et nombre de clients dans ces modèles. Cette dernière propriété implique l'existence d'un débit asymptotique non nul quand le nombre de files et le nombre de clients tendent vers l'infini.

**Keywords:** Equivalence, dualité client/serveur, réseaux de files d'attente en tandem fermés, serveur infini avec reséquencement, files d'attente avec blocage, réseau symétrique.

# 1 Introduction

In this paper, we establish equivalence relations between closed tandem queueing networks. We consider four types of models: single-server infinite-capacity-buffer queues, infinite-server queues with resequencing, single-server unit-capacity-buffer queues with blocking before service, and single-server unit-capacity-buffer queues with blocking after service.

Using a customer/server duality we show that in a closed tandem network consisting of single-server infinite-capacity-buffer queues, the case when service times are associated with customers, i.e. customer-dependent service, and the case when service times are associated with servers, i.e. server-dependent service, yield equivalent performance characteristics. This property allows one to analyze systems with customer-dependent service times from known results on systems with server-dependent service times.

We further show that for closed tandem queueing networks with customer-dependent service, a system consisting of single-server infinite-capacity-buffer queues (resp. infinite-server queues with resequencing) and a system consisting of single-server unit-capacity-buffer queues with blocking before service (resp. blocking after service) have equivalent performance behaviors. Since exact analysis and bounds are sometimes easier to obtain in one model than in another model, these equivalences allow us to analyze one system via the other one.

We then apply all these equivalence relations to symmetric closed tandem networks, where all the service times are independent and identically distributed random variables. We obtain various new performance analysis results. In particular, we obtain a closed-form expression for the throughput of networks with unit-capacity-buffer queues and blocking before service when the service times are exponentially distributed. We show the monotonicity of throughput (of queues) with respect to the number of queues and number of customers in these models. This last property in turn implies the existence of nonzero asymptotic throughput when the number of queues and number of customers go to infinity. We also prove that in a network with single-server infinite-capacity-buffer queues, the stationary cycle time distribution and the total throughput of the network are symmetric in the number of queues and the number of customers. Moreover, for any fixed sum of these two numbers, if the service times have a PERT type distribution (cf. Baccelli and Liu [6]), the total throughput of the network is concave in these numbers and is maximized when their absolute difference is smaller than or equal to 1.

The idea of interchanging the role between customers and servers was introduced before in Rego and Szpankowski [18] and Weber [21]. In [18], this idea was used to study polling systems with unit-capacity buffers by looking at a dual closed queueing network. Whereas in [21], this idea was used to show that the makespan is independent of the order of the queues in an open tandem network. We will remark that the customer/server duality is somewhat related to the customer/hole duality concept which was first introduced in Gordon and Newell [14] for closed tandem queues with finite-capacity buffers, and used thereafter by Akyildiz [1], Ammar and Gershwin [2] and Dallery, Liu and Towsley [11], etc., in related models.

The monotonicity of throughput (of queues) with respect to the number of queues and number of customers in these models are related to (and some of them are derived from) the monotonicity properties established by Baccelli and Liu [6]. Our proof of the existence of nonzero asymptotic throughput uses the large deviation bound of Baccelli and Konstantopoulos [3]. The symmetry of the stationary cycle time distribution with respect to the number of queues and the number of customers in a network with single-server infinite-capacity-buffer queues was first proved by Schassberger and Daduna [19] for exponential service times (they in fact obtained a closed form expression which is an Erlang distribution). The concavity result comes from our equivalence result and a concavity property in [6].

All models considered in this paper can be represented by stochastic strongly connected marked graphs (see Appendix A which illustrates how this can be done). The reader is referred to Baccelli, Cohen, Olsder and Quadrat [4] for the recent development on marked graphs. For networks with blocking, the reader is referred to Onvural [16] for a survey and Perros [17] for a bibliography. The model with resequencing buffers is an example of synchronization constraints. For an overview of results on queueing models for systems with synchronization constraints we refer to Baccelli and Makowski [7].

The organization of the paper is as follows. In Section 2, we describe in detail the models and notation used in the paper. In Section 3, we present the customer/server duality and establish equivalence relations between different models. In Section 4, we apply these properties to symmetric queueing models in order to obtain new quantitative and qualitative results. Finally, in Section 5, we provide further applications of our equivalence properties and extensions.

## 2 Models and Notation

In this paper, we consider closed tandem queueing networks where a number of customers, say  $M$ , pass a number of service stations, say  $N$ , in a cyclic order. In a closed tandem network with  $N$  stations and  $M$  customers, the stations are labeled  $1, 2, \dots, N$ , and the customers are labeled  $1, 2, \dots, M$ . The topology of the system (i.e. the order in which the stations are ordered) is described by the flow  $(n_1 \rightarrow n_2 \rightarrow \dots \rightarrow n_N \rightarrow n_1)$ , where  $n_j \in \{1, 2, \dots, N\}$ ,  $1 \leq j \leq N$ .

The following closed tandem queueing network models will be analyzed in the paper.

**Model 1 :** Single-server queues with infinite-capacity buffers.

In each station there is one single server and a buffer of infinite size, i.e. at least  $M$ . At each station, the customers are served in first come first serve (FCFS) order.

**Model 2 :** Infinite-server queues with resequencing.

Without loss of generality, the customers are labeled in the reverse direction of the network flow, i.e., for all  $2 \leq i \leq M$ , customer  $i$  is initially either in the same station as  $i - 1$  or in

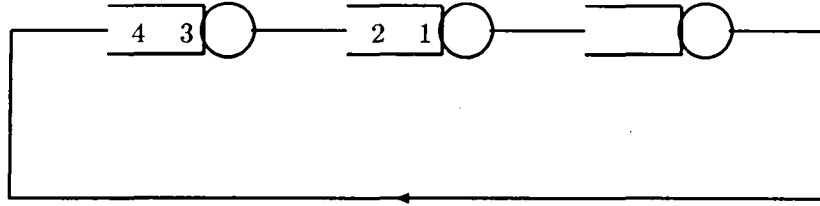


Figure 1: Model 1: single-server infinite-capacity-buffer queues.

an upstream station of that of  $i - 1$ . In each station there are infinitely many (i.e. at least  $M$ ) servers. After the  $n$ -th service completion, a customer, say  $i$ ,  $1 \leq i \leq M$ , can leave the station if and only if customers  $1, \dots, i - 1$  have left the same station for the  $n$ -th time, and customers  $i + 1, \dots, M$  have left the same station for the  $n - 1$ -st time,  $n \geq 1$ . In this way, customers do not overtake each other.

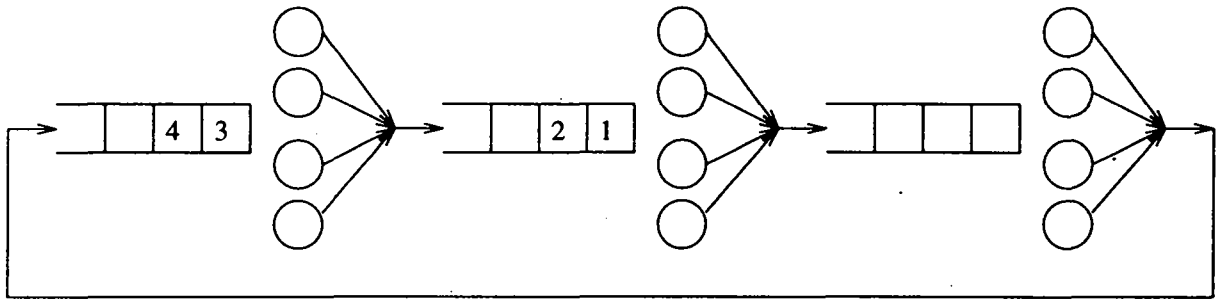


Figure 2: Model 2: infinite-server queues with resequencing.

**Model 3 :** Single-server queues with unit-capacity buffers and blocking before service.

In each station there is one single server and a waiting buffer of size 1. By convention, we assume that the servers have no buffer, i.e., the customer in service remains in the waiting buffer. Thus,  $N > M$ . The blocking mechanism is the so-called blocking before service (or communication blocking): a customer starts service at a station only when the downstream station is empty.

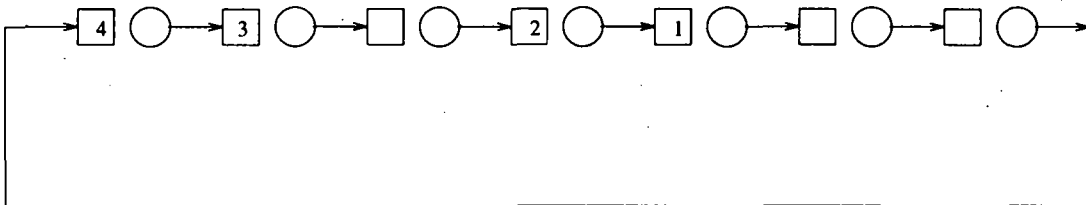


Figure 3: Models 3 and 4: single-server queues with unit-capacity buffers.



**Model 4 :** Single-server queues with unit-capacity buffers and blocking after service.

The model is similar to the previous one. In each station there is one single server and a waiting buffer of size 1. Again, the servers are assumed to have no buffer. Thus,  $N > M$ . The blocking mechanism is the so-called blocking after service (or manufacturing blocking): a customer starts service at a station as soon as it enters the station. At the completion of its service, it leaves the station only when the downstream station becomes empty. The server is blocked during the time interval between a service completion and the departure of the corresponding customer from the station.

All these models can be represented by stochastic (strongly connected) marked graphs (see Appendix A which illustrates how this can be done). Therefore, the results established in Baccelli and Liu [5] concerning the existence and the uniqueness of the throughput and stationary cycle times are readily applicable to our systems (see below). Moreover, the stochastic monotonicity and concavity properties obtained by Baccelli and Liu [6] and the large deviation bounds derived in Baccelli and Konstantopoulos [3] will be applied to our systems together with the equivalence properties established in the current paper to prove the existence of nonzero asymptotic throughput.

Let  $\sigma_{n,k}^s$  be the service duration of the  $k$ -th service initiated at station  $n$ ,  $1 \leq n \leq N$ ,  $k \geq 1$ . Let  $\sigma_{m,k}^c$  be the service duration of the  $k$ -th service of customer  $m$ ,  $1 \leq m \leq M$ ,  $k \geq 1$ . The system will be said to have server-dependent service if the service times are specified by the sequences  $\{\sigma_{n,k}^s\}_{k=1}^\infty$ ,  $1 \leq n \leq N$ , and to have customer-dependent service if the service times are specified by the sequences  $\{\sigma_{m,k}^c\}_{k=1}^\infty$ ,  $1 \leq m \leq M$ .

In this paper, the network will be called symmetric if all the service times are independent and identically distributed (i.i.d.) with the same (cumulated) distribution function. In a symmetric system there is no distinction between server-dependent and customer-dependent services.

In general, in a system with server-dependent service, the identity of the customers is not important. The topology of the system determines its performance characteristics. However, in a system with customer-dependent service, its performance characteristics are determined by the manner in which the customers are ordered and are spaced among them.

Therefore, for the models 1 and 2, i.e. closed tandem queueing networks with infinite-capacity-buffer queues, the state of the system ( $n_1 \rightarrow n_2 \rightarrow \dots \rightarrow n_N \rightarrow n_1$ ) at time  $t \geq 0$ , referred to as  $\mathcal{S}_t$ , will be represented by the following vector:

$$\left( m_{n_1, q_{n_1}}, \dots, m_{n_1, 1}, -n_1, m_{n_2, q_{n_2}}, \dots, m_{n_2, 1}, -n_2, \dots, m_{n_N, q_{n_N}}, \dots, m_{n_N, 1}, -n_N \right), \quad (2.1)$$

where  $q_n$  is the number of customers in station  $n$  (including those in service, if any), and  $m_{n,v}$  is the identity of the  $v$ -th customer in station  $n$  according to FCFS order,  $1 \leq n \leq N$ ,  $1 \leq v \leq q_n$ . Note that in the above representation, the negative integers represent the stations and the

positive integers the customers. For example, the state of the systems in Figure 1 and 2 is given by  $(4, 3, -1, 2, 1, -2, -3)$ .

By convention, we will assume that any rotational shift of the state yields an equivalent representation. For example, the state

$$(-n_N, m_{n_1, q_{n_1}}, \dots, m_{n_1, 1}, -n_1, m_{n_2, q_{n_2}}, \dots, m_{n_2, 1}, -n_2, \dots, m_{n_N, q_{n_N}}, \dots, m_{n_N, 1})$$

is considered the same as that of (2.1).

The state representation is said to be in canonical form if it is in the form of (2.1) and if  $n_N = N$ . Note that the canonical representation of a state is unique.

Two state representations  $\mathcal{S}$  and  $\mathcal{S}'$  are said to be equivalent, denoted by  $\mathcal{S} \equiv \mathcal{S}'$ , if they can be obtained one from the other by rotational shifts. In other words,  $\mathcal{S} \equiv \mathcal{S}'$  if  $\mathcal{S}$  and  $\mathcal{S}'$  have the same canonical representation.

For the models 3 and 4, i.e. closed tandem queueing networks with single-server unit-capacity-buffer queues, since there is at most one customer at a station, the state of the system at time  $t \geq 0$ , referred to as  $\mathcal{S}_t$ , will be represented by the following canonical form:

$$(m_{n_1}, m_{n_2}, \dots, m_{n_N}) \quad (2.2)$$

where  $n_N = N$ , and  $m_n$  is the identity of the customer at station  $n$  if any, and  $m_n = 0$  otherwise. For example, the state of the system in Figure 3 is given by  $(4, 3, 0, 2, 1, 0, 0)$ .

Let  $B_{n,k}^s$  (resp.  $C_{n,k}^s$ ) be the time epoch when the  $k$ -th service is initiated (resp. completed) at station  $n$ ,  $1 \leq n \leq N$ ,  $k \geq 1$ . Let  $B_{m,k}^c$  (resp.  $C_{m,k}^c$ ) be the time epoch when customer  $m$  starts (resp. completes) its  $k$ -th service,  $1 \leq m \leq M$ ,  $k \geq 1$ .

In closed tandem networks, the throughputs (defined as the number of service completions per unit of time) of the different stations are identical. Let  $\theta$  denote the throughput of a station, and  $\Theta$  the throughput of the network. Assume that the service times are integrable and that the sequence  $\{\sigma_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty$  or the sequence  $\{\sigma_{m,k}^c, 1 \leq m \leq M\}_{k=1}^\infty$  is stationary and ergodic. It was shown in Baccelli and Liu [5] that

$$\theta = \lim_{k \rightarrow \infty} \frac{k}{C_{n,k}^s} = \lim_{k \rightarrow \infty} \frac{k}{E[C_{n,k}^s]}, \quad a.s., \quad 1 \leq n \leq N. \quad (2.3)$$

where *a.s.* denotes “almost surely”. Similarly,

$$\theta = \lim_{k \rightarrow \infty} \frac{k}{C_{m,k}^c} = \lim_{k \rightarrow \infty} \frac{k}{E[C_{m,k}^c]}, \quad a.s., \quad 1 \leq m \leq M. \quad (2.4)$$

Throughout this paper we assume that these limits exist. It is clear that for a network of  $N$  stations,

$$\Theta = N \cdot \theta. \quad (2.5)$$

The cycle time in the network, defined from the viewpoint of a customer, is the sum of the successive sojourn times of a tagged customer at stations  $n_1, n_2, \dots, n_N$ . It can also be defined from the viewpoint of a station, i.e. the time between two successive departures of a tagged customer from a certain station. Conditions under which the cycle times converge to stationary random variables were obtained in [5]. In particular, if, in addition to the stationary and ergodic assumption on the service times, there is a station  $n$  such that its service times  $\{\sigma_{n,k}^s\}_{k=1}^\infty$  have infinite support and are independent of the other ones, then the cycle times converge (in coupling or in total variation) to unique stationary random variables. Throughout this paper, we will assume that the stationary cycle times exist, and are denoted by  $\phi_n^s$  (resp.  $\phi_m^c$ ) for station  $n$  (resp. customer  $m$ ). For models 1 and 3, they are defined by

$$\phi_n^s \stackrel{d}{=} \lim_{k \rightarrow \infty} (C_{n,k+N}^s - C_{n,k}^s), \quad 1 \leq n \leq N, \quad (2.6)$$

$$\phi_m^c \stackrel{d}{=} \lim_{k \rightarrow \infty} (C_{m,k+M}^c - C_{m,k}^c), \quad 1 \leq m \leq M, \quad (2.7)$$

whereas for models 2 and 4,

$$\phi_n^s \stackrel{d}{=} \lim_{k \rightarrow \infty} (B_{n,k+N}^s - B_{n,k}^s), \quad 1 \leq n \leq N, \quad (2.8)$$

$$\phi_m^c \stackrel{d}{=} \lim_{k \rightarrow \infty} (B_{m,k+M}^c - B_{m,k}^c), \quad 1 \leq m \leq M, \quad (2.9)$$

where  $\stackrel{d}{=}$  denotes equality in distribution.

In a symmetric system there is no distinction between the servers and between the customers. We will thus use  $\phi$  to denote the stationary cycle time.

It is easily seen that the first moments of these cycle times are the same in the same network. Let  $E[\phi]$  denote this mean cycle time. We have

$$\theta = \frac{M}{E[\phi]}, \quad \text{and} \quad \Theta = \frac{MN}{E[\phi]}. \quad (2.10)$$

In the sequel, whenever necessary, a super index  $i$ ,  $i = 1, 2, 3, 4$ , will be used in the above notation and quantities to indicate the type of models of the network. They will also be parameterized by the number of stations and the number of customers in the system. For instance,  $\mathcal{N}^1(N, M)$  denotes a model-1 network with  $N$  single-server infinite-capacity-buffer queues and  $M$  customers. The symbol  $C_{n,k}^{s,1}(N, M)$  denotes the  $k$ -th service completion time at queue  $n$  in such a network.

### 3 Duality and Equivalences

#### 3.1 Customer/Server Duality and Equivalence between Server-Dependent Service and Customer-Dependent Service

We first consider queueing networks of model 1, i.e. single-server queues with infinite-capacity buffers. We will show that in a certain sense, server-dependent service and customer-dependent service are equivalent. In order to do that, we introduce the notion of customer/server duality.

Let  $\mathcal{N}$  be a network of  $N$  stations and  $M$  customers, with topology  $(N \rightarrow N-1 \rightarrow \dots \rightarrow 1 \rightarrow N)$ , and the initial state  $\mathcal{S}_0$

$$\begin{aligned} \mathcal{S}_0 = & (q_1^\circ + \dots + q_N^\circ, \dots, q_1^\circ + \dots + q_{N-1}^\circ + 1, -N, \\ & q_1^\circ + \dots + q_{N-1}^\circ, \dots, q_1^\circ + \dots + q_{N-2}^\circ + 1, -N+1, \dots, q_1^\circ, \dots, 1, -1), \end{aligned} \quad (3.1)$$

where  $q_n^\circ$  is the initial number of customers in station  $n$ ,  $1 \leq n \leq N$ .

Let  $\widetilde{\mathcal{N}}$  be the customer/server dual (or simply, dual) of  $\mathcal{N}$ . The dual network  $\widetilde{\mathcal{N}}$  of  $\mathcal{N}$  is obtained from  $\mathcal{N}$  in the following way:

- reverse the flow of  $\mathcal{N}$ ;
- replace customer  $i$  of  $\mathcal{N}$  by station  $i$ ;
- replace station  $i$  of  $\mathcal{N}$  by customer  $i$ .

The resulting network  $\widetilde{\mathcal{N}}$ , viz., the dual of  $\mathcal{N}$ , has  $M$  stations and  $N$  customers. The customers of  $\widetilde{\mathcal{N}}$  in between two stations are considered as the customers of the downstream station (in the direction of flow of  $\widetilde{\mathcal{N}}$ ). Thus, the topology of  $\widetilde{\mathcal{N}}$  is  $(1 \rightarrow 2 \rightarrow \dots \rightarrow M \rightarrow 1)$ , and the initial state of  $\widetilde{\mathcal{N}}$  is obtained by reversing the string of  $\mathcal{S}_0$  and changing the signs of its components:

$$\begin{aligned} \widetilde{\mathcal{S}}_0 = & (1, -1, \dots, -q_1^\circ, \dots, N-1, -q_1^\circ - \dots - q_{N-2}^\circ - 1, \dots, -q_1^\circ - \dots - q_{N-1}^\circ, \\ & N, -q_1^\circ - \dots - q_{N-1}^\circ - 1, \dots, -q_1^\circ - \dots - q_N^\circ). \end{aligned} \quad (3.2)$$

Such a dual transformation is illustrated in Figure 4, where the white circles represent stations and the black ones the customers. In the example,  $N = 3$ ,  $M = 4$ . The initial state of the original system is  $(-3, 4, 3, -2, 2, 1, -1)$ , whereas the initial state of the dual system is  $(1, -1, -2, 2, -3, -4, 3)$ .

In the sequel, the variables associated with the dual network  $\widetilde{\mathcal{N}}$  will be denoted by the symbol tilde “ $\sim$ ”. For instance,  $\tilde{\sigma}_{m,k}^c$  is the service duration of the  $k$ -th service of customer  $m$

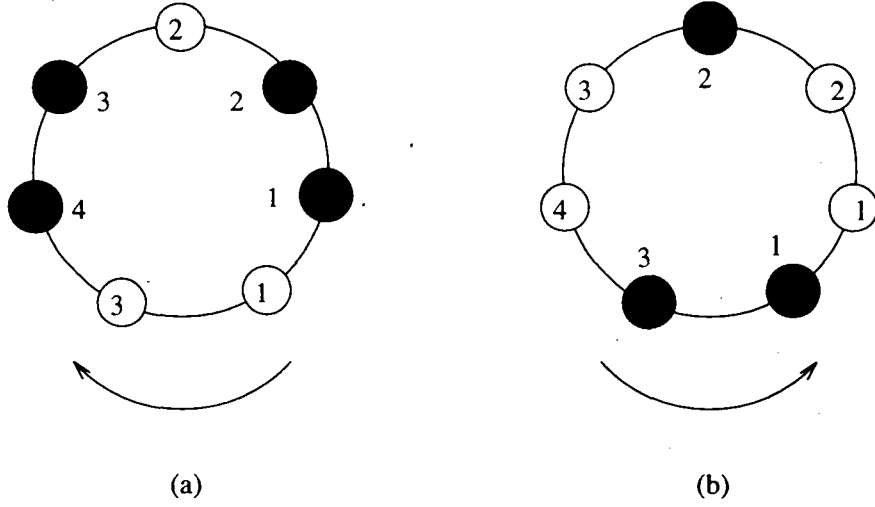


Figure 4: Customer/server dual. (a) Original network. (b) Dual network.

in the dual network  $\tilde{\mathcal{N}}$ . The operator “ $\mathbf{R}$ ” denotes the reversion of a string, and the operator “ $-$ ” changes the sign of the components in a string. For instance,

$$(1, -1, -2, 2, -3, -4, 3) = -\mathbf{R}(-3, 4, 3, -2, 2, 1, -1).$$

We are now in a position to prove the following equivalence relations between server-dependent service and customer-dependent service in closed tandem queueing networks with single server and infinite-capacity buffers.

**Theorem 3.1** *Let  $\mathcal{N}$  be a model-1 network with  $N$  stations and  $M$  customers, and  $\tilde{\mathcal{N}}$  be its (customer/server) dual.*

- If for all  $1 \leq n \leq N$ , and all  $k \geq 1$ ,  $\sigma_{n,k}^s = \tilde{\sigma}_{n,k}^c$ , then

$$S_t \equiv -\mathbf{R}(\tilde{S}_t), \quad t \geq 0, \quad (3.3)$$

$$B_{n,k}^s = \tilde{B}_{n,k}^c, \quad 1 \leq n \leq N, \quad k \geq 1, \quad (3.4)$$

$$C_{n,k}^s = \tilde{C}_{n,k}^c, \quad 1 \leq n \leq N, \quad k \geq 1. \quad (3.5)$$

- If for all  $1 \leq m \leq M$ , and all  $k \geq 1$ ,  $\sigma_{m,k}^c = \tilde{\sigma}_{m,k}^s$ , then

$$S_t \equiv -\mathbf{R}(\tilde{S}_t), \quad t \geq 0, \quad (3.6)$$

$$B_{m,k}^c = \tilde{B}_{m,k}^s, \quad 1 \leq m \leq M, \quad k \geq 1, \quad (3.7)$$

$$C_{m,k}^c = \tilde{C}_{m,k}^s, \quad 1 \leq m \leq M, \quad k \geq 1. \quad (3.8)$$

**Proof.** The assertions are shown by induction on the event (i.e. service commencement or completion) epochs  $0 = t_1 < t_2 < \dots < t_k < \dots$ . We only consider the relations (3.3–3.5). The proof of the others is analogous.

For  $t_1 = 0$ , it is readily checked from the definition of the customer/server dual network,

$$\mathcal{S}_0 \equiv -\mathbf{R}(\tilde{\mathcal{S}}_0).$$

Assume there is some integer  $v \geq 1$  such that relation (3.3) holds for all  $0 \leq t \leq t_v$ , and relation (3.4) holds for all  $n, k$  such that  $B_{n,k}^s < t_v$ .

In  $\mathcal{N}$  (resp.  $\tilde{\mathcal{N}}$ ), a server, say  $n$ , can start its  $k$ -th service by time  $t_v$  if and only if the number, say  $m$ , on the left of  $-n$  in  $\mathcal{S}_{t_v}$  (resp.  $\tilde{\mathcal{S}}_{t_v}$ ) (or one of its equivalent representations) is positive (representing the first customer to be served or being served). Since  $\tilde{\mathcal{N}}$  is obtained by reversing the flow of  $\mathcal{N}$  and by interchanging the servers and the customers in  $\mathcal{N}$ , the inductive assumption  $\mathcal{S}_{t_v} \equiv -\mathbf{R}(\tilde{\mathcal{S}}_{t_v})$  implies that customer  $m$  can start service at station  $n$  in  $\mathcal{N}$  at time  $t_v$  if and only if customer  $n$  can start service at station  $m$  in  $\tilde{\mathcal{N}}$  at time  $t_v$ . Thus, relation (3.4) holds for all  $n, k$  such that  $B_{n,k}^s \leq t_v$ . This in turn implies that relation (3.3) holds for all  $t_v \leq t \leq t_{v+1}$ , as for all  $1 \leq n \leq N$ , and all  $k \geq 1$ ,  $\sigma_{n,k}^s = \tilde{\sigma}_{n,k}^c$ .

Therefore, by induction, (3.3) holds for all  $t \geq 0$ , and (3.4) holds for all  $1 \leq n \leq N$  and all  $k \geq 1$ . Relation (3.5) is a consequence of (3.4).  $\blacksquare$

**Corollary 3.1** *Let  $\mathcal{N}$  be a model-1 network with  $N$  stations and  $M$  customers, and  $\tilde{\mathcal{N}}$  be its (customer/server) dual. Then, under the conditions of Theorem 3.1,*

$$\phi_n^s \stackrel{d}{=} \tilde{\phi}_n^c, \quad \phi_m^c \stackrel{d}{=} \tilde{\phi}_m^s, \quad 1 \leq n \leq N, \quad 1 \leq m \leq M, \quad (3.9)$$

and

$$N \cdot \theta = \Theta = \tilde{\Theta} = M \cdot \tilde{\theta}. \quad (3.10)$$

**Proof.** The cycle time in the original network defined from the viewpoint of customer  $m$ , i.e.  $C_{m,k+M}^c - C_{m,k}^c$  (resp. from the viewpoint of station  $n$ , i.e.  $C_{n,k+N}^s - C_{n,k}^s$ ) corresponds to the cycle time in the dual network defined from the viewpoint of station  $m$ , i.e.  $\tilde{C}_{m,k+M}^s - \tilde{C}_{m,k}^s$  (resp. from the viewpoint of customer  $n$ , i.e.  $\tilde{C}_{n,k+N}^c - \tilde{C}_{n,k}^c$ ). Relation (3.9) now follows from the facts that these transient cycle times (i.e. differences of service completion times) are identical according to Theorem 3.1, and that they converge to unique stationary random variables (cf. [5]).

Due to Theorem 3.1, at any time, the number of service completions in  $\mathcal{N}$  and in  $\widetilde{\mathcal{N}}$  are the same. Thus, relation (3.10) holds. ■

The customer/server duality concept enables us to analyze queueing models in which service times are related to customers instead of to queues. To illustrate this fact we present the following result.

**Theorem 3.2** *Let  $\mathcal{N}$  be a model-1 network with  $N$  queues and  $M$  customers. Assume that the service times that customer  $i$  demands at the different queues are i.i.d. with exponential distribution with parameter  $\mu_i$ ,  $i = 1, \dots, M$ . Assume further that all service times are mutually independent. Then, the steady-state cycle times of the different customers have the same distribution, whose Laplace-Stieltjes transform is given by*

$$f(s) = \sum_{(j_1, \dots, j_M) \in Z(M, N-1)} p(j_1, \dots, j_M) \prod_{i=1}^M \left( \frac{\mu_i}{\mu_i + s} \right)^{j_i+1}, \quad (3.11)$$

where

$$Z(M, N-1) = \left\{ (j_1, \dots, j_M) : j_i \geq 0, \sum_{i=1}^M j_i = N-1 \right\} \quad (3.12)$$

and

$$p(j_1, \dots, j_M) = \prod_{i=1}^M \mu_i^{-j_i} \left( \sum_{(j_1, \dots, j_M) \in Z(M, N-1)} \prod_{i=1}^M \mu_i^{-j_i} \right)^{-1}. \quad (3.13)$$

**Proof.** The result comes from relation (3.9) and the results of Schassberger and Daduna [19]. ■

Note that in Theorem 3.2, the cycle time is independent of the order in which the customers are served in a particular station (or the initial arrangement of the customers in the network). However, as remarked in Boxma [8], in a general model-1 network with server-dependent service, the cycle time distribution depends on the order in which the queues are visited unless all service times are exponential. Thus, in a general model-1 network with customer-dependent service, the cycle time distribution depends on the order in which the customers are served in a particular station.

### 3.2 Equivalences between a Model without Blocking and a Model with Blocking

We now establish some equivalence relations between networks with infinite-capacity buffers and those with unit-capacity buffers. Roughly speaking, the results below indicate that in a model-1 (resp. model-2) network, if the customers as well as the stations are replaced by single-server

unit-capacity-buffer queues (see Figure 3 where the network corresponds to such a transformation from the networks in Figures 1 and 2), then the resulting model-3 (resp. model-4) network has equivalent performance characteristics.

Throughout this subsection, we will consider networks with customer-dependent service. In such a case, the identities of the servers have no importance. Hence, for networks of models 3 and 4, we will use the following form of state representation (cf. (2.2))

$$(m_1, m_2, \dots, m_N),$$

where  $m_n$  is the identity of the customer at a station if any, and  $m_n = 0$  otherwise. It is assumed that  $m_{n+1}$  (resp.  $m_{n-1}$ ) is the customer (if any) at the downstream (resp. upstream) station of customer  $m_n$ . Note that the additions and subtractions on indices should be understood as modulo  $N$ , i.e.,  $m_{N+1}$  denotes  $m_1$  and  $m_0$  denotes  $m_N$ .

We will use the symbol “ $T$ ” to denote the operator that changes the negative components in a string to zero. For instance,

$$T(1, -1, -2, 2, -3, -4, 3) = (1, 0, 0, 2, 0, 0, 3).$$

**Theorem 3.3** *Let  $\mathcal{N}^1(N, M)$  be a model-1 network with  $N$  stations and  $M$  customers, and  $\mathcal{N}^3(N + M, M)$  a model-3 network with  $N + M$  stations and  $M$  customers. If  $T(\mathcal{S}_0^1(N, M)) \equiv \mathcal{S}_0^3(N + M, M)$ , and if for all  $1 \leq m \leq M$  and all  $k \geq 1$ ,  $\sigma_{m,k}^{c,1}(N, M) = \sigma_{m,k}^{c,3}(N + M, M)$ , then*

$$T(\mathcal{S}_t^1(N, M)) \equiv \mathcal{S}_t^3(N + M, M), \quad t \geq 0, \quad (3.14)$$

$$B_{m,k}^{c,1}(N, M) = B_{m,k}^{c,3}(N + M, M), \quad 1 \leq m \leq M, \quad k \geq 1, \quad (3.15)$$

$$C_{m,k}^{c,1}(N, M) = C_{m,k}^{c,3}(N + M, M), \quad 1 \leq m \leq M, \quad k \geq 1. \quad (3.16)$$

**Proof.** Note first that both  $\mathcal{S}_t^1(N, M)$  and  $\mathcal{S}_t^3(N + M, M)$  have  $N + M$  components.

The assertions are shown by induction on the event (i.e. service commencement or completion) epochs  $0 = t_1 < t_2 < \dots < t_k < \dots$ .

For  $t_1 = 0$ , (3.14) trivially holds. Assume there is some integer  $v \geq 1$  such that relation (3.14) holds for all  $0 \leq t \leq t_v$ , and relation (3.15) holds for all  $m, k$  such that  $B_{m,k}^{c,1}(N, M) < t_v$ .

In  $\mathcal{N}^1(N, M)$  (resp.  $\mathcal{N}^3(N + M, M)$ ), a customer, say  $m$ , can start its  $k$ -th service by time  $t_v$  if and only if the number, say  $n$ , on the right of  $m$  in  $\mathcal{S}_{t_v}^1(N, M)$  (resp.  $\mathcal{S}_{t_v}^3(N + M, M)$ ) (or one of its equivalent representations) is negative (resp. zero). In such a case, for  $\mathcal{N}^1(N, M)$ ,  $-n$  is the server on which customer  $m$  can start service, whereas for  $\mathcal{N}^3(N + M, M)$ ,  $n = 0$  means that the downstream queue of the station where customer  $m$  is located is empty so that customer



$m$  can start service (recall that the blocking mechanism in  $\mathcal{N}^3(N + M, M)$  is blocking before service). Therefore, the inductive assumption  $T(\mathcal{S}_{t_v}^1(N, M)) \equiv \mathcal{S}_{t_v}^3(N + M, M)$  implies that customer  $m$  can start service in  $\mathcal{N}^1(N, M)$  at time  $t_v$  if and only if customer  $m$  can start service in  $\mathcal{N}^3(N + M, M)$  at time  $t_v$ . Thus, relation (3.15) holds for all  $m, k$  such that  $B_{m,k}^{c,1}(N, M) \leq t_v$ . This in turn implies that relation (3.14) holds for all  $t_v \leq t \leq t_{v+1}$ , as for all  $1 \leq m \leq M$ , and all  $k \geq 1$ ,  $\sigma_{m,k}^{c,1}(N, M) = \sigma_{m,k}^{c,3}(N + M, M)$ .

Therefore, by induction, (3.14) holds for all  $t \geq 0$ , and (3.15) holds for all  $1 \leq n \leq N$  and all  $k \geq 1$ . Relation (3.16) is a consequence of (3.15). ■

**Corollary 3.2** *Let  $\mathcal{N}^1(N, M)$  be a model-1 network with  $N$  stations and  $M$  customers, and  $\mathcal{N}^3(N + M, M)$  a model-3 network with  $N + M$  stations and  $M$  customers. If  $T(\mathcal{S}_0^1(N, M)) \equiv \mathcal{S}_0^3(N + M, M)$ , and if for all  $1 \leq m \leq M$  and all  $k \geq 1$ ,  $\sigma_{m,k}^{c,1}(N, M) = \sigma_{m,k}^{c,3}(N + M, M)$ , then*

$$N \cdot \theta^1(N, M) = \Theta^1(N, M) = \Theta^3(N + M, M) = (N + M) \cdot \theta^3(N + M, M). \quad (3.17)$$

Combining the above equivalence properties (Theorem 3.3 and Corollary 3.2) and the duality properties (Theorem 3.1 and Corollary 3.1), we can analyze model-3 networks with customer-dependent service by analyzing model-1 networks with server-dependent service. In particular, when the service times are exponential, we obtain product-form solutions, as is illustrated in the following theorem.

**Theorem 3.4** *Let  $\mathcal{N}^3(N, M)$  be a model-3 network with  $N$  queues and  $M$  customers,  $N > M$ . Assume that the service times are mutually independent, and that the service times of customer  $i$  are i.i.d. with exponential distribution of parameter  $\mu_i$ ,  $1 \leq i \leq M$ . Then the network throughput of  $\mathcal{N}^3(N, M)$  is given by*

$$\Theta^3(N, M) = M \cdot \frac{G_M(N - M - 1)}{G_M(N - M)}, \quad (3.18)$$

where  $G_M(K)$  is the normalizing constant of the product-form solution of the model-1 network with  $M$  stations and  $K$  customers and server-dependent exponential service times:

$$G_M(K) = \sum_{\substack{n_1 + \dots + n_M = K; \\ n_1, \dots, n_M \geq 0}} \left( \prod_{i=1}^M \frac{1}{\mu_i^{n_i}} \right). \quad (3.19)$$

**Proof.** Let  $\mathcal{N}^1(N - M, M)$  be a model-1 network with  $M$  customers and  $N - M$  single-server infinite-capacity-buffer queues. Assume that the service times in  $\mathcal{N}^1(N - M, M)$  are mutually

independent, and that the service times of customer  $i$  are i.i.d. with exponential distribution with parameter  $\mu_i$ ,  $1 \leq i \leq M$ . Let  $\widetilde{\mathcal{N}}^1(M, N - M)$  be the dual of  $\mathcal{N}^1(N - M, M)$ , with  $M$  stations and  $N - M$  customers. Then the service times in  $\widetilde{\mathcal{N}}^1(M, N - M)$  are mutually independent, and the service times of station  $i$  are i.i.d. with exponential distribution of parameter  $\mu_i$ ,  $1 \leq i \leq M$ .

By coupling the service times in these networks (see the proof of Theorem 4.1 below), we obtain from Corollaries 3.1 and 3.2 that the network throughputs of these systems are identical:

$$\Theta^3(N, M) = \Theta^1(N - M, M) = \widetilde{\Theta}^1(M, N - M).$$

Since  $\widetilde{\mathcal{N}}^1(M, N - M)$  has a product-form solution, we obtain (see e.g. Gelenbe and Mitrani [13, p. 103]) for the throughput of a station of  $\widetilde{\mathcal{N}}^1(M, N - M)$ ,

$$\widetilde{\theta}^1(M, N - M) = \frac{G_M(N - M - 1)}{G_M(N - M)}.$$

Thus

$$\Theta = \widetilde{\Theta}^1(M, N - M) = M \cdot \frac{G_M(N - M - 1)}{G_M(N - M)}.$$

■

Note that the computation of the normalizing constant  $G_M(K)$  can be carried out efficiently in  $O(MK)$  steps (see Buzen [9]).

Similar to the equivalence between model-1 and model-3 networks, the model-2 (infinite servers with resequencing) and model-4 (unit-capacity buffers with blocking after service) networks are equivalent:

**Theorem 3.5** *Let  $\mathcal{N}^2(N, M)$  be a model-2 network with  $N$  stations and  $M$  customers, and  $\mathcal{N}^4(N + M, M)$  a model-4 network with  $N + M$  stations and  $M$  customers. If  $T(\mathcal{S}_0^2(N, M)) \equiv \mathcal{S}_0^4(N + M, M)$ , and if for all  $1 \leq m \leq M$  and all  $k \geq 1$ ,  $\sigma_{m,k}^{c,2}(N, M) = \sigma_{m,k}^{c,4}(N + M, M)$ , then*

$$T(\mathcal{S}_t^2(N, M)) \equiv \mathcal{S}_t^4(N + M, M), \quad t \geq 0, \quad (3.20)$$

$$B_{m,k}^{c,2}(N, M) = B_{m,k}^{c,4}(N + M, M), \quad 1 \leq m \leq M, \quad k \geq 1, \quad (3.21)$$

$$C_{m,k}^{c,2}(N, M) = C_{m,k}^{c,4}(N + M, M), \quad 1 \leq m \leq M, \quad k \geq 1. \quad (3.22)$$

**Proof.** The proof is similar to that of Theorem 3.3, and is carried out by induction on the event (i.e. service commencement or completion) epochs. Note that in  $\mathcal{N}^2(N, M)$  and  $\mathcal{N}^4(N + M, M)$ , a departure time of a customer coincides with the service commencement time of the customer in the downstream station. Note also that in  $\mathcal{N}^2(N, M)$  and  $\mathcal{N}^4(N + M, M)$ , the  $k$ -th departure of

a customer, say  $m$ , occurs only after its service completion and the  $k$ -th departures of customers  $1, \dots, m-1$ . The detailed proof is left to the interested reader. ■

**Corollary 3.3** *Let  $\mathcal{N}^2(N, M)$  be a model-2 network with  $N$  stations and  $M$  customers, and  $\mathcal{N}^4(N+M, M)$  a model-4 network with  $N+M$  stations and  $M$  customers. If  $T(\mathcal{S}_0^2(N, M)) \equiv \mathcal{S}_0^4(N+M, M)$ , and if for all  $1 \leq m \leq M$  and all  $k \geq 1$ ,  $\sigma_{m,k}^{c,2}(N, M) = \sigma_{m,k}^{c,4}(N+M, M)$ , then*

$$N \cdot \theta^2(N, M) = \Theta^2(N, M) = \Theta^4(N+M, M) = (N+M) \cdot \theta^4(N+M, M). \quad (3.23)$$

Using the above equivalence properties, we can easily prove the following bounds on the throughput of model-4 networks.

**Theorem 3.6** *Let  $\mathcal{N}$  be a model-4 network with  $M$  customers and  $N$  single-server unit-capacity-buffer queues,  $N > M$ . Assume that the sequences of service times  $\{\sigma_{m,k}^c\}_{k=1}^\infty$ ,  $1 \leq m \leq M$ , are mutually independent and are stationary and ergodic. Then,*

$$\frac{M}{E \left[ \max_{1 \leq m \leq M} \sigma_{m,1}^c \right]} \leq \Theta \leq \sum_{m=1}^M \frac{1}{E[\sigma_{m,1}^c]}. \quad (3.24)$$

**Proof.** Let  $\mathcal{N}'$  be a model-2 network with  $M$  customers and  $N-M$  infinite-server queues with resequencing. The sequences of service times in  $\mathcal{N}'$  are coupled with those of  $\mathcal{N}$  in the sense that  $\sigma_{m,k}^c$  is also the service time of the  $k$ -th service of customer  $m$  in  $\mathcal{N}'$ . It then follows from Corollary 3.3 that  $\mathcal{N}$  and  $\mathcal{N}'$  have the same network throughput:

$$\Theta = \Theta'. \quad (3.25)$$

Let  $\mathcal{N}'_1$  and  $\mathcal{N}'_2$  be two model-2 networks having the same topology and the same service times as those of  $\mathcal{N}'$ . They differ from  $\mathcal{N}'$  only in that the resequencing mechanism is relaxed and enhanced, respectively. In  $\mathcal{N}'_1$ , there is no resequencing, whereas in  $\mathcal{N}'_2$ , the customer departures from any station are synchronized: a customer can leave station  $n$  for the  $k$ -th time only when all the customers have completed the  $k$ -th service at station  $n$ . It is easily seen (cf. Baccelli and Liu [6]) that  $\mathcal{N}'_1$  has a larger network throughput and  $\mathcal{N}'_2$  has a smaller network throughput than  $\mathcal{N}'$ :

$$\Theta'_2 \leq \Theta' \leq \Theta'_1. \quad (3.26)$$

Moreover, in  $\mathcal{N}'_1$ , the customers are always in service, so that  $\Theta'_1$  is the sum of the number of service completions of each customer per unit of time. Hence,

$$\Theta'_1 = \sum_{m=1}^M \frac{1}{E[\sigma_{m,1}^c]}. \quad (3.27)$$

In  $\mathcal{N}'_2$ , however, the customer arrivals to (as well as the departures from) any station are synchronized, so that the mean cycle time  $E[\phi'_2]$  is

$$E[\phi'_2] = E \left[ \sum_{k=1}^{N-M} \max_{1 \leq m \leq M} \sigma_{m,k}^c \right] = (N-M) E \left[ \max_{1 \leq m \leq M} \sigma_{m,1}^c \right].$$

Thus,

$$\Theta'_2 = \frac{M(N-M)}{E[\phi'_2]} = \frac{M}{E \left[ \max_{1 \leq m \leq M} \sigma_{m,1}^c \right]}. \quad (3.28)$$

Combining relations (3.25—3.28) yields the desired result. ■

## 4 Applications to Symmetric Networks

In this section, we present the applications of the above duality and equivalence properties to symmetric networks, i.e. networks where all the service times are i.i.d. with the same distribution function. In such networks, there is no distinction between server-dependent and customer-dependent services. We are interested in performance measures such as throughput and cycle time in these networks.

### 4.1 Symmetry, Monotonicity and Concavity Properties

**Theorem 4.1** *Let  $\mathcal{N}$  be a symmetric model-1 network with  $N$  stations and  $M$  customers. Then the distribution of the stationary cycle time and the network throughput are symmetric in  $M$  and  $N$ .*

**Proof.** Let  $\widetilde{\mathcal{N}}$  be the dual of  $\mathcal{N}$  with i.i.d. service times. Assume that the distribution functions of the service times in  $\mathcal{N}$  and  $\widetilde{\mathcal{N}}$  are the same. We can thus couple the service times of  $\mathcal{N}$  and  $\widetilde{\mathcal{N}}$  in such a way that for all  $1 \leq n \leq N$ , and all  $k \geq 1$ ,  $\sigma_{n,k}^s = \tilde{\sigma}_{n,k}^c$ . (Note that more precisely, we should say that there exists a probability space such that for all  $1 \leq n \leq N$ , and all  $k \geq 1$ ,  $\sigma_{n,k}^s = \tilde{\sigma}_{n,k}^c$  a.s.) It then follows from Corollary 3.1 that  $\phi_n^s$  (resp.  $\phi_m^c$ ) and  $\tilde{\phi}_n^c$  (resp.  $\tilde{\phi}_m^s$ ) are identical in law,  $1 \leq n \leq N$ ,  $1 \leq m \leq M$ .

Note that the only difference between  $\mathcal{N}$  and  $\widetilde{\mathcal{N}}$  is that the number of servers and the number of customers are interchanged (the direction of flow has no importance as the network is symmetric). Therefore, we conclude that the distribution of  $\phi$  is symmetric in  $M$  and  $N$ . As a consequence, the network throughput of a symmetric model-1 network is also symmetric in  $M$

and  $N$  (cf. (2.10)). ■

**Remark:** As mentioned in the introduction, the result of Theorem 4.1 is well-known for exponential service times (cf. Schassberger and Daduna [19]). In that case they even showed that the cycle time distribution depends on  $N$  and  $M$  only through the sum of them.

We show below that the station throughput of a symmetric closed tandem queueing network is monotonically decreasing in its number of stations and number of customers under some conditions.

**Theorem 4.2** *Let  $\mathcal{N}^i(N_i, M_i)$  (resp.  $\widehat{\mathcal{N}}^i(\widehat{N}_i, \widehat{M}_i)$ ) be a symmetric model- $i$  network with  $N_i$  stations and  $M_i$  customers (resp.  $\widehat{N}_i$  stations and  $\widehat{M}_i$  customers),  $i = 1, 2, 3, 4$ . Assume that for each  $i = 1, 2, 3, 4$ , the service times in  $\mathcal{N}^i(N_i, M_i)$  and  $\widehat{\mathcal{N}}^i(\widehat{N}_i, \widehat{M}_i)$  are equivalent in law. If for  $i = 1, 2, 3, 4$ , conditions (4.1), (4.2), (4.3), (4.4) hold respectively,*

$$\widehat{N}_1 \geq N_1, \quad \widehat{M}_1 \geq M_1, \quad \widehat{N}_1 - N_1 \geq \widehat{M}_1 - M_1; \quad (4.1)$$

$$\widehat{N}_2 \geq N_2, \quad \widehat{M}_2 \geq M_2; \quad (4.2)$$

$$\widehat{N}_3 \geq N_3, \quad \widehat{M}_3 \geq M_3, \quad \widehat{N}_3 - N_3 \geq 2(\widehat{M}_3 - M_3), \quad \frac{\widehat{N}_3}{\widehat{M}_3} \leq \frac{N_3}{M_3}; \quad (4.3)$$

$$\widehat{N}_4 \geq N_4, \quad \widehat{M}_4 \geq M_4, \quad \widehat{N}_4 - N_4 \geq \widehat{M}_4 - M_4, \quad \frac{\widehat{N}_4}{\widehat{M}_4} \leq \frac{N_4}{M_4}, \quad (4.4)$$

then the station throughput of  $\widehat{\mathcal{N}}^i(\widehat{N}_i, \widehat{M}_i)$  is smaller than that of  $\mathcal{N}^i(N_i, M_i)$ :

$$\hat{\theta}^i(\widehat{N}_i, \widehat{M}_i) \leq \theta^i(N_i, M_i), \quad i = 1, 2, 3, 4. \quad (4.5)$$

**Proof.** The cases  $i = 1$  and  $i = 2$  are special cases of Theorems 5.1 and 5.2 below, respectively.

For case  $i = 3$ , let  $\mathcal{N}'(N_3 - M_3, M_3)$  (resp.  $\widehat{\mathcal{N}}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3)$ ) be a symmetric model-1 network with  $N_3 - M_3$  stations and  $M_3$  customers (resp.  $\widehat{N}_3 - \widehat{M}_3$  stations and  $\widehat{M}_3$  customers). The service times in  $\mathcal{N}^3(N_3, M_3)$ ,  $\mathcal{N}'(N_3 - M_3, M_3)$ ,  $\widehat{\mathcal{N}}^3(\widehat{N}_3, \widehat{M}_3)$  and  $\widehat{\mathcal{N}}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3)$  are equivalent in law. It then follows from Corollary 3.2 that

$$\begin{aligned} \theta^3(N_3, M_3) &= \frac{N_3 - M_3}{N_3} \cdot \theta'(N_3 - M_3, M_3), \\ \hat{\theta}^3(\widehat{N}_3, \widehat{M}_3) &= \frac{\widehat{N}_3 - \widehat{M}_3}{\widehat{N}_3} \cdot \hat{\theta}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3). \end{aligned}$$

Under condition (4.3),

$$\widehat{N}_3 - \widehat{M}_3 \geq N_3 - M_3, \quad \widehat{M}_3 \geq M_3, \quad (\widehat{N}_3 - \widehat{M}_3) - (N_3 - M_3) \geq \widehat{M}_3 - M_3.$$

Thus, applying (4.5) (for case  $i = 1$ ) to networks  $\mathcal{N}'(N_3 - M_3, M_3)$  and  $\widehat{\mathcal{N}}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3)$  implies that

$$\hat{\theta}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3) \leq \theta'(N_3 - M_3, M_3).$$

It then follows that

$$\begin{aligned} \hat{\theta}^3(\widehat{N}_3, \widehat{M}_3) &= \frac{\widehat{N}_3 - \widehat{M}_3}{\widehat{N}_3} \cdot \hat{\theta}'(\widehat{N}_3 - \widehat{M}_3, \widehat{M}_3) \\ &\leq \frac{\widehat{N}_3 - \widehat{M}_3}{\widehat{N}_3} \cdot \theta'(N_3 - M_3, M_3) \\ &\leq \frac{N_3 - M_3}{N_3} \cdot \theta'(N_3 - M_3, M_3) \\ &= \theta^3(N_3, M_3), \end{aligned}$$

where the last inequality comes from the fact that  $\widehat{N}_3/\widehat{M}_3 \leq N_3/M_3$  according to (4.3).

The case  $i = 4$  can be shown in an analogous way in using Corollary 3.3 and the relation (4.5) (for case  $i = 2$ ). The detailed proof is omitted.  $\blacksquare$

**Remark:** For model-1 and model-2 networks, the monotonicity properties are established under more general statistical assumptions in Theorems 5.1 and 5.2 in the next section. An application of the monotonicity properties of Theorem 4.2 will be presented later on in the paper to show the nonzero asymptotic station throughput when the number of stations and the number of customers tend to infinity.

According to a result of Baccelli and Liu [6], we know that in model-1 networks, when the service times have a PERT type distribution, the station throughput and the network throughput are concave in the number of customers in the system. A distribution function is of PERT type if it is a distribution function of the length of the critical path of a PERT graph where the weights are random variables with exponential distributions. In particular, the class of PERT type distributions includes the exponential and Erlang distributions. Note that a concavity property holds for nonsymmetric model-1 networks with server-dependent service, i.e., the servers can have different service time distributions.

We now use the equivalence property (Corollary 3.2) to show the following concavity property.

**Theorem 4.3** *Let  $\mathcal{N}$  be a symmetric model-1 network with  $N$  stations and  $M$  customers. Let  $L = M + N$  be fixed. Assume that the service time distribution is of PERT type. Then, the network throughput of  $\mathcal{N}$  is concave in  $M$  and  $N$ . Furthermore, the network throughput is maximized when  $|M - N| \leq 1$ .*

**Proof.** Let  $\mathcal{N}'$  be a symmetric model-3 network with  $M$  customers and  $L = N + M$  queues. The service times in  $\mathcal{N}'$  have the same distribution as those in  $\mathcal{N}$ . It follows from Corollary 3.2

that the networks  $\mathcal{N}$  and  $\mathcal{N}'$  have the same network throughput:  $\Theta = \Theta'$ .

Applying further the result of Baccelli and Liu [6] implies that  $\Theta'$  is concave in  $M$ , so is  $\Theta$ . Moreover,  $\Theta$  is concave in  $N = L - M$ .

Since  $\Theta$  is symmetric in  $M$  and  $N = L - M$  according to Theorem 4.1, we obtain that  $\Theta$  is maximized when  $M = \lfloor L/2 \rfloor$  or  $M = \lceil L/2 \rceil$ , where  $\lfloor x \rfloor$  is the integer part of  $x$ , and  $\lceil x \rceil$  is the smallest integer greater than or equal to  $x$ . In other words,  $M = N$  or  $M = N \pm 1$ . ■

**Remark:** For model-3 networks, the concavity and the maximization of the network throughput were obtained in [11] for nonsymmetric server-dependent service.

## 4.2 Exact Analysis and Bounds of Throughput

We now consider the throughput of queueing networks with blocking. We first obtain the following closed-form solution of the throughput in a symmetric closed tandem network with single-server unit-capacity-buffer queues and exponential service times.

**Theorem 4.4** *Let  $\mathcal{N}$  be a symmetric model-3 network with  $M$  customers and  $N$  single-server unit-capacity-buffer queues,  $N > M$ . Assume that all the service times are i.i.d. with exponential distribution of parameter  $\mu$ . Then the network throughput  $\Theta$  and the station throughput  $\theta$  of  $\mathcal{N}$  are given by*

$$\Theta = \mu \cdot \frac{M(N - M)}{N - 1}, \quad \theta = \mu \cdot \frac{M(N - M)}{N(N - 1)}. \quad (4.6)$$

**Proof.** It follows from Theorem 3.4 that

$$\Theta = M \cdot \frac{G_M(N - M - 1)}{G_M(N - M)}, \quad (4.7)$$

where  $G_M(K)$  is the normalizing constant of the product-form solution of the model-1 network with  $M$  stations and  $K$  customers with exponential service of parameter  $\mu$ . Thus (cf. (3.19)),

$$G_M(K) = \sum_{\substack{n_1 + \dots + n_M = K; \\ n_1, \dots, n_M \geq 0}} \mu^{-K}. \quad (4.8)$$

There are  $\binom{M + K - 1}{M - 1}$  terms in the right-hand side of the above equation. To see this, we note that a vector  $(n_1, \dots, n_M)$  corresponds to the binary chain

$$\underbrace{1 \dots 1}_{n_1} 0 \underbrace{1 \dots 1}_{n_2} 0 \dots 0 \underbrace{1 \dots 1}_{n_M}$$

with  $K$  ones and  $M - 1$  zeros. The state vector  $(n_1, \dots, n_M)$  is thus determined by the positions of the zeros among  $K + M - 1$  letters.

Thus, it follows from (4.7) that

$$\Theta = M \cdot \frac{G_M(N - M - 1)}{G_M(N - M)} = M \cdot \frac{\binom{N - 2}{M - 1} \cdot \mu^{-N+M+1}}{\binom{N - 1}{M - 1} \cdot \mu^{-N+M}} = \mu \cdot \frac{M(N - M)}{N - 1}.$$

■

Consider now the equivalence properties between model-2 (infinite server with resequencing) and model-4 (blocking after service) networks. We assume here that these networks are symmetric. In the literature, there is no analytical method for these networks, even under the assumption of exponential service times. Various approximate solutions were proposed for the model-4 networks, see [16] for a survey. According to the equivalence properties (Theorem 3.5 and Corollary 3.3), these approximate solutions are applicable to the model-2 networks.

Moreover, using these equivalence properties, the bounds obtained in one model can be transmitted to another. The bounds in Theorem 3.6 illustrated such idea. We present here some more examples.

**Corollary 4.1** *Let  $\mathcal{N}$  be a symmetric model-4 network with  $M$  customers and  $N$  queues,  $N > M$ . Assume all the service times are i.i.d. exponentially distributed with parameter  $\mu$ . Then,*

$$\max \left\{ \mu \cdot \frac{M(N - M)}{N - 1}, \mu \cdot \frac{M}{\sum_{i=1}^M \frac{1}{i}} \right\} \leq \Theta \leq M\mu. \quad (4.9)$$

**Proof.** First, as a consequence of Theorem 3.6,

$$\mu \cdot \frac{M}{\sum_{i=1}^M \frac{1}{i}} \leq \Theta \leq M\mu. \quad (4.10)$$

It was shown in [12] that the blocking-after-service mechanism yields a larger throughput than the corresponding network with blocking-before-service mechanism. Thus, for the symmetric model-4 network  $\mathcal{N}$  with  $M$  customers,  $N$  queues and exponential service times, a lower bound of the throughput can be obtained from the corresponding model-3 network using Theorem 4.4:

$$\Theta \geq \mu \cdot \frac{M(N - M)}{N - 1}. \quad (4.11)$$

Combining (4.10) and (4.11) yields (4.9). ■



Using the lower bound in (4.9) together with Corollary 3.3, we obtain

**Corollary 4.2** *Let  $\mathcal{N}$  be a symmetric model-2 network with  $M$  customers and  $N$  queues. Assume that all the service times are i.i.d. exponentially distributed with parameter  $\mu$ . Then,*

$$\Theta \geq \max \left\{ \mu \cdot \frac{M(N-M)}{N-1}, \mu \cdot \frac{M}{\sum_{i=1}^M \frac{1}{i}} \right\}. \quad (4.12)$$

### 4.3 Asymptotic Throughput

We now consider the asymptotic station throughputs of the closed tandem networks when the number of stations and the number of customers tend to infinity. Our results rely on the monotonicity properties of Theorem 4.2 and the large deviation bound established in Baccelli and Konstantopoulos [3] (see also [4, pp. 394-404]).

Let  $\sigma$  be a random variable. Assume that the moment generating function of  $\sigma$  exists for some  $z_0 > 0$ ,

$$h(z) \stackrel{\text{def}}{=} E(e^{z\sigma}) < \infty, \quad \forall z \leq z_0. \quad (4.13)$$

Let  $H(x)$  be the Cramer-Legendre transform for the distribution function of  $\sigma$ ,

$$H(x) \stackrel{\text{def}}{=} \inf_{z \in \mathbb{R}} (\log(h(z)) - zx). \quad (4.14)$$

Let

$$\gamma \stackrel{\text{def}}{=} \inf \{x \mid x > E[\sigma], H(x) + \log(2) < 0\} \quad (4.15)$$

**Lemma 4.1** *For  $i = 1, 2, 3, 4$ , let  $\mathcal{N}^i(N_i, M_i)$  be a symmetric model- $i$  network whose service times have the same law as the random variable  $\sigma$ . If  $N_i = M_i$  for  $i = 1, 2$ , and  $N_i = 2M_i$  for  $i = 3, 4$ , then the station throughput of  $\mathcal{N}^i(N_i, M_i)$  is bounded below by  $1/\gamma$ :*

$$\theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{1}{\gamma}, \quad i = 1, 2, 3, 4. \quad (4.16)$$

**Proof.** It has been shown in [11] that in a strongly connected marked graph, two reachable initial markings yield the same throughput. Thus, the throughputs of the networks  $\mathcal{N}^i(N_i, M_i)$ ,  $i = 1, 2, 3, 4$ , are independent of the initial queue lengths. We can therefore assume, without loss of generality, that in  $\mathcal{N}^1(N_1, M_1)$ , at time zero, there is one and only one customer in each station.

As is illustrated in Appendix A (Figure 6), the model-1 networks can be represented by a subclass of stochastic Petri nets. Note that for sake of simplicity we assume that any transition

has no simultaneous firing. This corresponds to the Petri nets with “recycled” transitions. Therefore, the out-degree of a Petri net representing the model-1 network is 2 (including the one exit arc for the “recycling place”).

By applying now the large deviation bound of [3], we obtain that

$$\theta^1(\mathcal{N}^1(M_1, M_1)) \geq \frac{1}{\gamma}.$$

As a consequence of the equivalence property between model-1 and model-3 networks (cf. Corollary 3.2), we get that

$$\theta^3(\mathcal{N}^3(2M_3, M_3)) = \theta^1(\mathcal{N}^1(M_3, M_3)) \geq \frac{1}{\gamma}.$$

As we mentioned in the previous subsection, a model-2 (resp. model-4) network has larger throughput than the model-1 (resp. model-3) network with the same number of stations and the same number of customers. Therefore,

$$\begin{aligned} \theta^2(\mathcal{N}^2(M_2, M_2)) &\geq \theta^1(\mathcal{N}^1(M_2, M_2)) \geq \frac{1}{\gamma}, \\ \theta^4(\mathcal{N}^4(2M_4, M_4)) &\geq \theta^3(\mathcal{N}^3(2M_4, M_4)) \geq \frac{1}{\gamma}. \end{aligned}$$

■

We are now in a position to establish the following existence property of asymptotic throughput.

**Theorem 4.5** *For  $i = 1, 2, 3, 4$ , let  $\mathcal{N}^i(N_i, M_i)$  be a symmetric model- $i$  network whose service times have the same law as the random variable  $\sigma$ . Let  $q \geq 1$  be a positive integer constant. If  $M_i = M$ ,  $N_i = qM$  for  $i = 1, 2$ , and  $N_i = (q + 1)M$  for  $i = 3, 4$ , then the asymptotic station throughput of  $\mathcal{N}^i(N_i, M_i)$  when  $M$  goes to infinity exists, denoted by  $\bar{\theta}^i(q)$ , and is bounded below by  $1/(q\gamma)$ :*

$$\bar{\theta}^i(q) = \lim_{M \rightarrow \infty} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{1}{q\gamma}, \quad i = 1, 2, 3, 4. \quad (4.17)$$

*If, further, the random variable  $\sigma$  has a PERT type distribution, then*

$$\bar{\theta}^i(q) = \lim_{M \rightarrow \infty} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}, \quad i = 1, 2, 3, 4. \quad (4.18)$$

**Proof.** Owing to Theorem 4.2, we have

$$\begin{aligned} \theta^i(\mathcal{N}^i(q(M+1), M+1)) &\leq \theta^i(\mathcal{N}^i(qM, M)), \quad i = 1, 2, \quad M \geq 1, \\ \theta^i(\mathcal{N}^i((q+1)(M+1), M+1)) &\leq \theta^i(\mathcal{N}^i((q+1)M, M)), \quad i = 3, 4, \quad M \geq 1. \end{aligned}$$

These monotonicities imply the existence of the asymptotic throughputs

$$\bar{\theta}^i(q) = \lim_{M \rightarrow \infty} \theta^i(\mathcal{N}^i(qM, M)), \quad i = 1, 2, \quad (4.19)$$

$$\bar{\theta}^i(q) = \lim_{M \rightarrow \infty} \theta^i(\mathcal{N}^i((q+1)M, M)), \quad i = 3, 4. \quad (4.20)$$

According to Lemma 4.1,  $\theta^1(\mathcal{N}^1(M, M)) \geq 1/\gamma$ , so that  $\bar{\theta}^1(1) \geq 1/\gamma$ .

For general  $q \geq 1$ , consider the network  $\mathcal{N}^1(qM, M)$ . Assume that the stations are labeled in a cyclic order so that the topology of the network is  $(1 \rightarrow 2 \rightarrow \dots \rightarrow qM)$ . Since the throughput is independent of the initial queue lengths, we can assume without loss of generality that at time zero, the customers are equally distributed in the sense that stations  $M, 2M, \dots, qM$  have one customer each, and the others have no customer. Let  $\tau = \sigma_1 + \dots + \sigma_q$ , where  $\sigma_1, \dots, \sigma_q$  are independent versions of random variable  $\sigma$ . Let

$$h_q(z) \stackrel{\text{def}}{=} E(e^{z\tau}) = (E(e^{z\sigma}))^q = (h(z))^q, \quad (4.21)$$

$$H_q(x) \stackrel{\text{def}}{=} \inf_{z \in \mathbf{R}} (\log(h_q(z)) - zx) = \inf_{z \in \mathbf{R}} (q \log(h(z)) - zx) = qH(x/q), \quad (4.22)$$

$$\gamma_q \stackrel{\text{def}}{=} \inf\{x \mid x > E[\tau], \quad H_q(x) + \log(2) < 0\}. \quad (4.23)$$

A simple calculation yields

$$\begin{aligned} q\gamma &= \inf\{qx \mid x > E[\sigma], \quad H(x) + \log(2) < 0\} \\ &= \inf\{y \mid y > qE[\sigma], \quad H(y/q) + \log(2) < 0\} \\ &= \inf\{y \mid y > E[\tau], \quad qH(y/q) + q\log(2) < 0\} \\ &\geq \inf\{y \mid y > E[\tau], \quad qH(y/q) + \log(2) < 0\} \\ &= \gamma_q. \end{aligned}$$

Using again the result of [3] implies that

$$\theta^1(\mathcal{N}^1(qM, M)) \geq \frac{1}{\gamma_q} \geq \frac{1}{q\gamma}. \quad (4.24)$$

Therefore,

$$\theta^2(\mathcal{N}^2(qM, M)) \geq \theta^1(\mathcal{N}^1(qM, M)) \geq \frac{1}{\gamma_q} \geq \frac{1}{q\gamma}, \quad (4.25)$$

$$\theta^3(\mathcal{N}^3((q+1)M, M)) = \theta^1(\mathcal{N}^1(qM, M)) \geq \frac{1}{\gamma_q} \geq \frac{1}{q\gamma}, \quad (4.26)$$

$$\theta^4(\mathcal{N}^4((q+1)M, M)) \geq \theta^3(\mathcal{N}^3((q+1)M, M)) \geq \frac{1}{\gamma_q} \geq \frac{1}{q\gamma}. \quad (4.27)$$

Equations (4.24)–(4.27), together with (4.19) and (4.20), readily imply relation (4.17).

If  $\sigma$  has a PERT type distribution, then, according to the concavity property of Theorem 4.3, we obtain that

$$\theta^1(\mathcal{N}^1(2qM, 2M)) \geq \frac{2}{q+1} \theta^1(\mathcal{N}^1((q+1)M, (q+1)M)) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}. \quad (4.28)$$

Relation (4.18) now follows as an immediate consequence. ■

**Corollary 4.3** *For  $i = 1, 2, 3, 4$ , let  $\mathcal{N}^i(N_i, M_i)$  be a symmetric model- $i$  network whose service times have the same law as the random variable  $\sigma$ . Let  $q$  be a positive integer constant such that  $q \geq 2$ . Then*

$$\bar{\theta}^i(q) = \lim_{\substack{N_i \rightarrow \infty, \quad M_i \rightarrow \infty \\ qM_i \leq N_i < (q+1)M_i}} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{1}{q\gamma}, \quad i = 1, 2, \quad (4.29)$$

$$\bar{\theta}^i(q) = \lim_{\substack{N_i \rightarrow \infty, \quad M_i \rightarrow \infty \\ (q+1)M_i \leq N_i < (q+2)M_i}} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{1}{q\gamma}, \quad i = 3, 4. \quad (4.30)$$

If, further, the random variable  $\sigma$  has a PERT type distribution, then

$$\bar{\theta}^i(q) = \lim_{\substack{N_i \rightarrow \infty, \quad M_i \rightarrow \infty \\ qM_i \leq N_i < (q+1)M_i}} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}, \quad i = 1, 2, \quad (4.31)$$

$$\bar{\theta}^i(q) = \lim_{\substack{N_i \rightarrow \infty, \quad M_i \rightarrow \infty \\ (q+1)M_i \leq N_i < (q+2)M_i}} \theta^i(\mathcal{N}^i(N_i, M_i)) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}, \quad i = 3, 4. \quad (4.32)$$

**Proof.** We only consider the case  $i = 1$ , i.e. the model-1 network  $\mathcal{N}^1(N_1, M_1)$ . The case  $i = 2$  is analogous and is omitted. The cases  $i = 3$  and  $i = 4$  are consequences of the cases 1 and 2, respectively, in view of Corollaries 3.2 and 3.3.

Owing to Theorem 4.2, we have that for all  $M_1 \geq 1$  and  $qM_1 \leq N_1 < (q+1)M_1$ ,

$$\theta^1(\mathcal{N}^1(2qM_1, 2M_1)) \leq \theta^1(\mathcal{N}^1(N_1, M_1)) \leq \theta^1(\mathcal{N}^1(qM_1, M_1)).$$

Using Theorem 4.5 and taking the limit when  $M_1$  and  $N_1$  go to infinity in the above inequalities imply that

$$\bar{\theta}^1(q) = \lim_{M_1 \rightarrow \infty} \theta^1(\mathcal{N}^1(2qM_1, 2M_1))$$

$$\begin{aligned}
&\leq \lim_{\substack{N_1 \rightarrow \infty, \quad M_1 \rightarrow \infty \\ qM_1 \leq N_1 < (q+1)M_1}} \theta^1(\mathcal{N}^1(N_1, M_1)) \\
&\leq \lim_{M_1 \rightarrow \infty} \theta^1(\mathcal{N}^1(qM_1, M_1)) \\
&= \bar{\theta}^1(q),
\end{aligned}$$

so that

$$\lim_{\substack{N_1 \rightarrow \infty, \quad M_1 \rightarrow \infty \\ qM_1 \leq N_1 < (q+1)M_1}} \theta^1(\mathcal{N}^1(N_1, M_1)) = \bar{\theta}^1(q).$$

Thus, relation (4.31) holds for  $i = 1$ . The proof is thus completed.  $\blacksquare$

## 5 Concluding Remarks, Extensions and Further Applications

In this paper, we have obtained duality and equivalence properties in closed tandem queueing networks. We have presented the customer/server duality and have shown the equivalence between customer-dependent and server-dependent services in model-1 networks. We have also shown the equivalence between model-1 and model-3 networks, and the equivalence between model-2 and model-4 networks.

These equivalence results hold in fact for nonzero initial workload, i.e., the customers and the servers have some work to finish at time zero before starting a complete service. Indeed, in such a case,  $\sigma_{n,1}^s$  (resp.  $\sigma_{m,1}^c$ ) can be considered as the remaining work of server  $n$  (resp. customer  $m$ ) at time zero. Thus, Theorems 3.1, 3.3 and 3.5 are still valid. Since the throughput and the stationary cycle time do not depend on the initial workload in general (see [5]), Corollaries 3.1, 3.2 and 3.3 remain true.

The customer/server duality presented in this paper is in certain sense related to the customer/hole duality introduced by Gordon and Newell [14]. It was shown by Dallery, Liu and Towsley [11] that an assembly/disassembly network with finite buffers and blocking before service (which is more general than the model-3 networks) is equivalent to its customer/hole dual where the roles of customers and holes (i.e. unoccupied space in the buffers) are interchanged and the flows are reversed. Using this result together with the equivalence between model-1 and model-3 networks established in this paper, we can show the customer/server duality for symmetric model-1 networks. In order to see this, we consider the network throughput  $\Theta^1(N, M)$  of a symmetric model-1 network  $\mathcal{N}^1(N, M)$  and the network throughput  $\tilde{\Theta}^1(M, N)$  of its customer/server dual  $\tilde{\mathcal{N}}^1(M, N)$ . According to Corollary 3.2,  $\Theta^1(N, M) = \Theta^3(N + M, M)$ . Applying the customer/hole duality to the symmetric model-3 network  $\mathcal{N}^3(N + M, M)$  implies  $\Theta^3(N + M, M) = \Theta^3(N + M, N)$ , where we note that the number of holes in  $\mathcal{N}^3(N + M, M)$

is  $N$ . Using again Corollary 3.2 we get  $\Theta^3(N + M, N) = \tilde{\Theta}^1(M, N)$ . Figure 5 illustrates these relations.

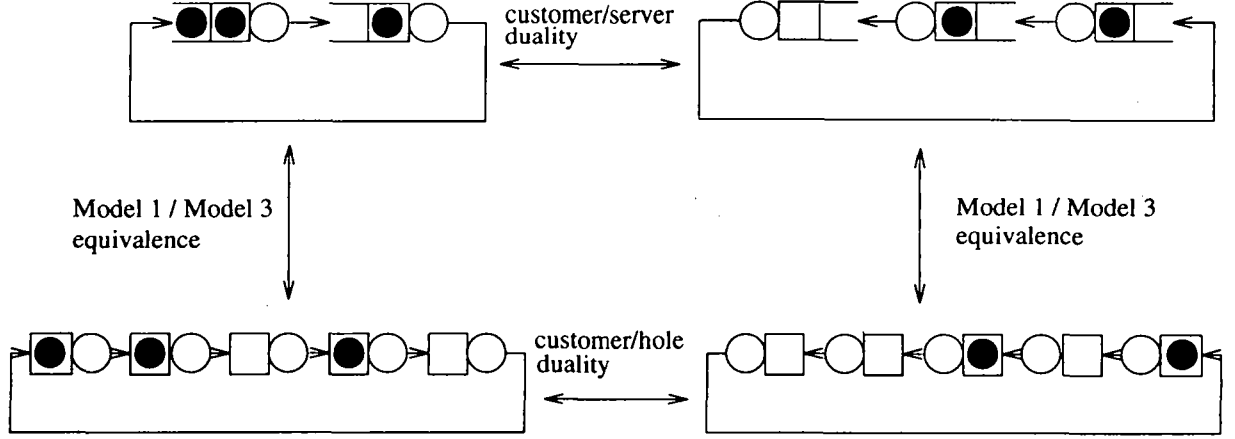


Figure 5: Relations between customer/server duality and customer/hole duality.

Some of the results obtained in Section 4 can be generalized using the notion of stochastic ordering (in particular the increasing convex ordering  $\leq_{\text{icx}}$ , see the definition below) and the stochastic comparison results obtained in [6].

A random vector  $X \in \mathbb{R}^n$  is said to be smaller than a random vector  $Y \in \mathbb{R}^n$  in the sense of increasing convex ordering, denoted by  $X \leq_{\text{icx}} Y$ , if for all increasing and convex functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $E[f(X)] \leq E[f(Y)]$ , provided the expectations exist. A sequence of random variables  $\{x_k\}_{k=1}^\infty$  is smaller than the sequence of random variables  $\{y_k\}_{k=1}^\infty$  in the sense of increasing convex ordering, denoted by  $\{x_k\}_{k=1}^\infty \leq_{\text{icx}} \{y_k\}_{k=1}^\infty$ , if for all  $n \geq 1$ ,  $\{x_k\}_{k=1}^n \leq_{\text{icx}} \{y_k\}_{k=1}^n$ .

For example, it is readily shown that a random variable which, with probability  $p$ , is a constant  $1/\mu$ , and with probability  $1 - p$ , is an exponentially distributed with parameter  $\mu$ , is smaller, in the sense of increasing convex ordering, than the random variable with exponential distribution of parameter  $\mu$ .

Thus, an application of Theorem 3.4 and the monotonicity property with respect to timing established in [6] implies

**Corollary 5.1** *Let  $\mathcal{N}$  be a model-3 network with  $M$  customers and  $N$  single-server unit-capacity-buffer queues,  $N > M$ . Assume that the service times are mutually independent, and that the service times of customer  $i$  are i.i.d. with mean  $1/\mu_i$ ,  $1 \leq i \leq M$ . Customer  $i$  has constant service time  $1/\mu_i$  with probability  $p_i$  and exponential service time of parameter  $\mu_i$  with probability  $1 - p_i$ ,  $0 \leq p_i \leq 1$ ,  $1 \leq i \leq M$ . Then the network throughput  $\Theta$  of  $\mathcal{N}$  is bounded below by*

$$\Theta \geq M \cdot \frac{G_M(N - M - 1)}{G_M(N - M)}, \quad (5.1)$$

where  $G_M(K)$  is determined by (3.19).

Concerning the asymptotic throughput, one easily observes, in view of the monotonicity property with respect to timing established in [6], that Theorem 4.5 and Corollary 4.3 remain valid if the service times are smaller than  $\sigma$  in the sense of increasing convex ordering. Moreover, owing to the following monotonicity properties (Theorems 5.1 and 5.2) of station throughput with respect to number of stations and number of customers in model-1 and model-2 networks, the results of Theorem 4.5 and Corollary 4.3 can be extended to asymmetric model-1 and model-2 networks (see Corollary 5.2 below).

In general, the network throughput and the station throughput of model-1 networks are monotonically increasing in its number of customers according to the monotonicity property with respect to initial markings obtained in [6]. However, the network throughput and the station throughput of a closed tandem queueing network (of one of the four models) are usually not monotone with respect to the number of stations. We prove in Appendix B that the following monotonicity of the station throughput in model-1 networks holds with respect to both the number of stations and the number of customers.

**Theorem 5.1** *Let  $\mathcal{N}(N, M)$  (resp.  $\widehat{\mathcal{N}}(\widehat{N}, \widehat{M})$ ) be a model-1 network with  $N$  stations and  $M$  customers (resp.  $\widehat{N}$  stations and  $\widehat{M}$  customers) and server-dependent service  $\{\sigma_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty$  (resp.  $\{\hat{\sigma}_{n,k}^s, 1 \leq n \leq \widehat{N}\}_{k=1}^\infty$ ). Assume that*

$$\widehat{N} \geq N, \quad \widehat{M} \geq M, \quad \widehat{N} - N \geq \widehat{M} - M.$$

*Assume further that the sequences  $\{\sigma_{n,k}^s\}_{k=1}^\infty, 1 \leq n \leq N$ , (resp.  $\{\hat{\sigma}_{n,k}^s\}_{k=1}^\infty, 1 \leq n \leq \widehat{N}$ ), are mutually independent, and that for all  $N < n \leq \widehat{N}$ , the upstream station of  $n$ , denoted  $n'$ ,  $1 \leq n' \leq \widehat{N}$ , is such that*

$$\{\sigma_{n',k}^s\}_{k=1}^\infty \leq_{\text{icx}} \{\hat{\sigma}_{n,k}^s\}_{k=1}^\infty.$$

*Then the station throughput of  $\widehat{\mathcal{N}}(\widehat{N}, \widehat{M})$  is smaller than that of  $\mathcal{N}(N, M)$ :*

$$\hat{\theta}(\widehat{N}, \widehat{M}) \leq \theta(N, M).$$

Note that the above monotonicity property is different from the monotonicity property with respect to topology [6]. Indeed, for two model-1 networks, one with more customers and more stations than the other, the corresponding marked graph of the former cannot, in general, be obtained by adding transitions and tokens in the corresponding marked graph of the latter. Thus, the result of [6] does not apply. Nevertheless, according to the representation of model-2 networks by stochastic marked graphs described in Appendix A, such a result does apply to model-2 networks, so that we can prove (see Appendix B) that in model-2 networks, the station throughput is monotonically decreasing in the number of stations and the number of customers.

**Theorem 5.2** Let  $\mathcal{N}(N, M)$  (resp.  $\widehat{\mathcal{N}}(\widehat{N}, \widehat{M})$ ) be a model-2 network with  $N$  stations and  $M$  customers (resp.  $\widehat{N}$  stations and  $\widehat{M}$  customers). Let  $\{\sigma_{m,n,k}, 1 \leq m \leq M, 1 \leq n \leq N\}_{k=1}^{\infty}$  (resp.  $\{\hat{\sigma}_{m,n,k}, 1 \leq m \leq \widehat{M}, 1 \leq n \leq \widehat{N}\}_{k=1}^{\infty}$ ) be the sequence of service times of customer  $m$  at station  $n$ . Assume that  $\widehat{N} \geq N$  and  $\widehat{M} \geq M$ . Assume further that

$$\{\sigma_{m,n,k}, 1 \leq m \leq M, 1 \leq n \leq N\}_{k=1}^{\infty} \leq_{\text{icx}} \{\hat{\sigma}_{m,n,k}, 1 \leq m \leq M, 1 \leq n \leq N\}_{k=1}^{\infty}.$$

Then the station throughput of  $\widehat{\mathcal{N}}(\widehat{N}, \widehat{M})$  is smaller than that of  $\mathcal{N}(N, M)$ :

$$\hat{\theta}(\widehat{N}, \widehat{M}) \leq \theta(N, M).$$

As a consequence, in model-1 and model-2 networks with server-dependent services, if the service times are increasing in the sense of increasing convex ordering, then the station throughput decreases when the number of stations and the number of customers increase. Hence, the results of Theorem 4.5 and Corollary 4.3 can be extended to asymmetric model-1 and model-2 networks.

**Corollary 5.2** For  $i = 1, 2$ , let  $\mathcal{N}^i(N_i, M_i)$  be an asymmetric model- $i$  network with server-dependent service times  $\{\sigma_{n,k}^{s,i}, 1 \leq n \leq N_i\}_{k=1}^{\infty}$ . Assume that the topology of the network is  $(1 \rightarrow 2 \rightarrow \dots \rightarrow N_i \rightarrow 1)$ . Assume also that the sequences of service times of the servers  $\{\sigma_{n,k}^{s,i}\}_{k=1}^{\infty}$  are mutually independent, and that each of which consists of i.i.d. random variables. If for all  $n \geq 1$  and  $i = 1, 2$ ,

$$\sigma_{n,1}^{s,i} \leq_{\text{icx}} \sigma_{n+1,1}^{s,i} \leq_{\text{icx}} \sigma,$$

then

$$\bar{\theta}^i(1) = \lim_{M \rightarrow \infty} \theta^i(\mathcal{N}^i(M, M)) \geq \frac{1}{\gamma}, \quad i = 1, 2, \quad (5.2)$$

and for any positive integer constant such that  $q \geq 2$ ,

$$\bar{\theta}^i(q) = \lim_{\substack{N \rightarrow \infty, \quad M \rightarrow \infty \\ qM \leq N < (q+1)M}} \theta^i(\mathcal{N}^i(N, M)) \geq \frac{1}{q\gamma}, \quad i = 1, 2, \quad (5.3)$$

where  $\sigma$  and  $\gamma$  satisfy (4.13), (4.14) and (4.15). If, further, the random variable  $\sigma$  has a PERT type distribution, then

$$\bar{\theta}^i(q) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}, \quad i = 1, 2. \quad (5.4)$$

For general closed tandem networks with buffers of arbitrary sizes, due to the monotonicity property with respect to initial markings and to timing [6], the throughput decreases when the buffer sizes decreases and/or when the service times increases in the sense of  $\leq_{\text{icx}}$  ordering. Therefore, Theorem 4.5 and Corollary 4.3 imply



**Corollary 5.3** *Let  $\mathcal{N}(N, M)$  be an arbitrary closed tandem queueing network with  $N$  single-server finite-capacity queues and  $M$  customers, and server-dependent service times  $\{\sigma_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty$ . The sizes of the buffers are arbitrary, and the blocking mechanism of any server can be either blocking before service or blocking after service. Assume that the sequences of service times of the servers  $\{\sigma_{n,k}^s\}_{k=1}^\infty$  are mutually independent, and that each of which consists of i.i.d. random variables. If for all  $n \geq 1$ ,  $\sigma_{n,1}^s \leq_{\text{icx}} \sigma$ , then*

$$\liminf_{M \rightarrow \infty} \theta(\mathcal{N}(2M, M)) \geq \frac{1}{\gamma}, \quad (5.5)$$

and for any positive integer constant such that  $q \geq 2$ ,

$$\liminf_{\substack{N \rightarrow \infty, \quad M \rightarrow \infty \\ (q+1)M \leq N < (q+2)M}} \theta(\mathcal{N}(N, M)) \geq \frac{1}{q\gamma}. \quad (5.6)$$

If, further, the random variable  $\sigma$  has a PERT type distribution, then

$$\liminf_{\substack{N \rightarrow \infty, \quad M \rightarrow \infty \\ (q+1)M \leq N < (q+2)M}} \theta(\mathcal{N}(N, M)) \geq \frac{2}{q+1} \cdot \frac{1}{\gamma}. \quad (5.7)$$

In a similar way, we can prove the existence of nonzero asymptotic station throughput for open tandem queueing networks with blocking. Let  $\mathcal{N}(N)$  be a tandem network consisting of  $N$  single-server queues with blocking. The first queue has an infinite-capacity buffer with infinite supply of customers so that the server is never starved. The last queue has an infinite-capacity downstream buffer so that the server is never blocked. The intermediate buffers have buffers with possibly finite size. The blocking mechanism of a server, when it has a finite-capacity downstream buffer, is either blocking before service or blocking after service. According to the monotonicity property with respect to topology [6], the throughput decreases when a new single-server queue (with finite or infinite buffer) is added at the end of the tandem network  $\mathcal{N}$ , the station throughput decreases. Therefore, asymptotic station throughput of such an open tandem queueing network  $\mathcal{N}(N)$  exists. Moreover, since for any given  $N$ , the throughput does not depend on the initial queue lengths of the queues, we can assume that there is one customer at each station initially. Therefore, using the large deviation bound of [3] implies that the station throughput is greater than  $1/\gamma$  when the service times are smaller than  $\sigma$  in the sense of  $\leq_{\text{icx}}$ . Hence,

**Corollary 5.4** *Let  $\mathcal{N}(N)$  be a tandem network consisting of  $N$  single-server queues with blocking. The first queue has an infinite-capacity buffer with infinite supply of customers and the last queue has infinite-capacity downstream buffer. The blocking mechanism of a server, when it has a finite-capacity downstream buffer, is either blocking before service or blocking after service.*

Assume that the sequences of service times of the servers  $\{\sigma_{n,k}^s\}_{k=1}^\infty$  are mutually independent, and that each of which consists of i.i.d. random variables. If for all  $n \geq 1$ ,  $\sigma_{n,1}^s \leq_{\text{icx}} \sigma$ , then

$$\bar{\theta}(\mathcal{N}) = \lim_{N \rightarrow \infty} \theta(\mathcal{N}(N)) \geq \frac{1}{\gamma}. \quad (5.8)$$

The above result proves a conjecture presented in [10] claiming that the throughput of a transfer line decreases to a nonzero value when the number of machines increases.

Note that in these asymptotic analyses of (station) throughput, we require that service times are bounded in  $\leq_{\text{icx}}$  sense by a random variable for which the moment generating function (4.13) exists for some  $z_0 > 0$ . A necessary and sufficient condition for the existence of the moment generating function of  $\sigma$  is that the tail distribution of  $\sigma$  is bounded by an exponential:

$$\exists \mu > 0, \exists x_0 > 0, \forall x > x_0: \quad P\{\sigma > x\} \leq e^{-\mu x}. \quad (5.9)$$

As we saw in the previous sections, for a symmetric model-1 network  $\mathcal{N}(N, M)$  with  $N$  stations and  $M$  customers, and with exponential service times of parameter  $\mu$ , its station throughput  $\theta^1(N, M)$  is expressed as

$$\theta^1(N, M) = \frac{G_N(M-1)}{G_N(M)},$$

where  $G_N(K)$  was defined in (4.8). Hence,

$$\theta^1(N, M) = \frac{G_N(M-1)}{G_N(M)} = \frac{\binom{N+M-2}{N-1} \mu^{-M+1}}{\binom{N+M-1}{N-1} \mu^{-M}} = \mu \cdot \frac{M}{N+M-1}. \quad (5.10)$$

When  $N = M$ ,  $\theta^1(N, N) \geq \mu/2$ . Therefore, when  $\sigma$  is smaller than an exponential random variable of parameter  $\mu$  in  $\leq_{\text{icx}}$  sense, we obtain another lower bound for the asymptotic station throughput: For any positive integer  $q \geq 1$ ,

$$\bar{\theta}^i(q) \geq \mu/2, \quad i = 1, 2, 3, 4. \quad (5.11)$$

## A Modeling Closed Tandem Queueing Networks by Stochastic Strongly Connected Marked Graphs

In this appendix, we illustrate how the queueing networks analyzed in this paper can be represented by stochastic marked graphs (which form a subset of stochastic Petri nets). Roughly speaking, a marked graph is a bipartite graph with a set of transitions and a set of places. The in-degree (resp. out-degree) of any place is at most one. The evolution of the graph is characterized by the circulation of tokens, which stay in places, and are consumed and produced by transitions. A transition is enabled to fire when there is at least one token in each of its upstream places. The firing consumes one token of each of these places and produces, after some firing time, one token into each of its downstream places. At any time, at most one firing is permitted at each transition.

The marked graph in Figure 6 represents the queueing network of Figure 1, where the transitions represent the servers, the places the buffers and the tokens the customers.

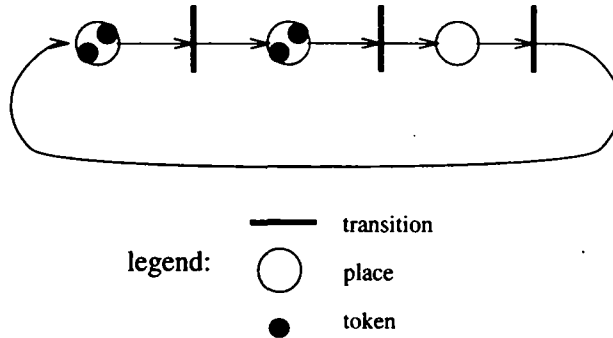


Figure 6: Stochastic marked graph corresponding to the model-1 network of Figure 1.

The marked graph in Figure 7 represents the queueing network of Figure 2, where the “thick-bar” transitions represent the servers, and the “thin-bar” transitions represent the departures and have zero firing time. Transition with label  $ij$  represents the server in station  $i$  dedicated to customer  $j$ ,  $1 \leq i \leq 3$ ,  $1 \leq j \leq 4$ .

The marked graph in Figure 8 represents the model-3 queueing network of Figure 3, where the transitions represent the servers, the places as well as the tokens in the middle row represent the waiting buffers and the customers, and the other places as well as the tokens represent the availability of the downstream waiting buffers.

The marked graph in Figure 9 represents the model-4 queueing network of Figure 3, where the “thick-bar” transitions represent the servers, and the “thin-bar” transitions represent the departures and have zero firing time. The places as well as the tokens in the middle row represent the waiting buffers and the customers, whereas the other places and the tokens represent the availability of the servers and the downstream waiting buffers.

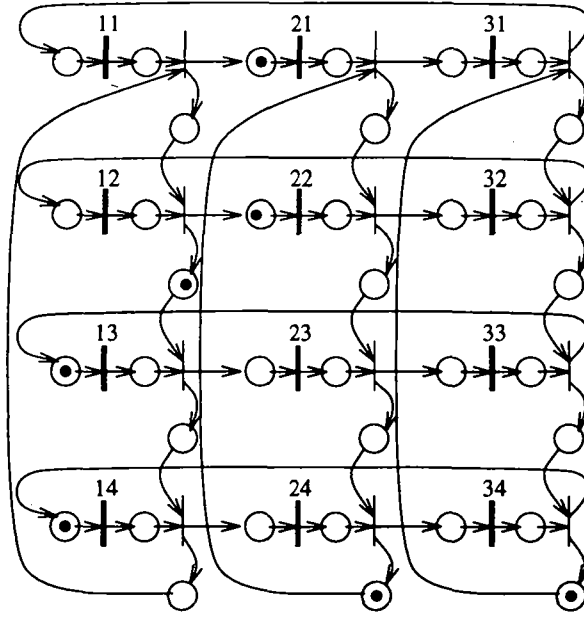


Figure 7: Stochastic marked graph corresponding to the model-2 network of Figure 2.

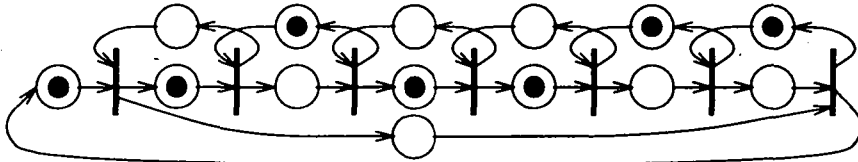


Figure 8: Stochastic marked graph corresponding to the model-3 network of Figure 3.

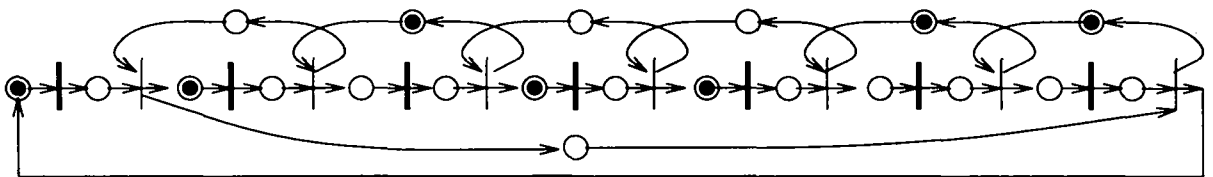


Figure 9: Stochastic marked graph corresponding to the model-4 network of Figure 3.

## B Monotonicity of Throughput in Number of Stations and Customers

The proofs of the monotonicity results are based on Strassen's Theorem [20] concerning increasing convex ordering, on the monotonicity results of stochastic marked graphs [6], and on the above correspondence between closed tandem queueing networks and stochastic marked graphs.

### Proof of Theorem 5.1.

Observe first that when a station with zero initial queue length is added in  $\mathcal{N}(N, M)$ , then the station throughput decreases, i.e.,

$$\theta^1(N+1, M) \leq \theta^1(N, M).$$

Such a relation is implied by the property of throughput monotonicity with respect to topology established in [6] for stochastic marked graphs. Indeed, the stochastic marked graph corresponding to  $\mathcal{N}(N+1, M)$  can be obtained by adding transitions and places to the stochastic marked graph corresponding to  $\mathcal{N}(N, M)$ . For example, Figure 10 contains a model-1 network obtained by inserting a station in between the two left most stations in Figure 1. By adding one transition and two places (Figure 11-(a)) in the marked graph (Figure 6) corresponding to the network of Figure 1, and by further removing a redundant place of the resulting marked graph (Figure 11-(a)), we obtain the marked graph (Figure 11-(b)) corresponding to the network of Figure 10.

Therefore, in what follows, we will only consider the case that  $\hat{N} - N = \hat{M} - M$ . Moreover, we will only consider the case that  $\hat{N} - N = \hat{M} - M = 1$ , as the general case can be completed by a simple induction. Without loss of generality, we assume that the topology of  $\mathcal{N}(N, M)$  is  $(1 \rightarrow 2 \rightarrow \dots \rightarrow N \rightarrow 1)$ , and that of  $\hat{\mathcal{N}}(N+1, M+1)$  is  $(1 \rightarrow 2 \rightarrow \dots \rightarrow N+1 \rightarrow 1)$ . Thus, the assumption of the theorem indicates that

$$\{\sigma_{N,k}^s\}_{k=1}^\infty \leq_{\text{icx}} \{\hat{\sigma}_{N+1,k}^s\}_{k=1}^\infty.$$

Owing to the property of throughput monotonicity with respect to timing [6] in stochastic marked graphs, we can further assume, without loss of generality, that the service times of stations  $i$  in the two networks  $\mathcal{N}(N, M)$  and  $\hat{\mathcal{N}}(N+1, M+1)$  have the same law for  $i = 1, \dots, N$ . Under the mutual independence assumption of the service times, we know that there exists a probability space such that

$$\{\sigma_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty = \{\hat{\sigma}_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty, \quad a.s. \quad (\text{B.1})$$

component-wise. Moreover, according to Strassen's Theorem [20] on increasing convex ordering, we obtain that

$$\{\sigma_{N,k}^s\}_{k=1}^\infty \leq E \left[ \{\hat{\sigma}_{N+1,k}^s\}_{k=1}^\infty \mid \{\sigma_{N,k}^s\}_{k=1}^\infty \right], \quad a.s. \quad (\text{B.2})$$

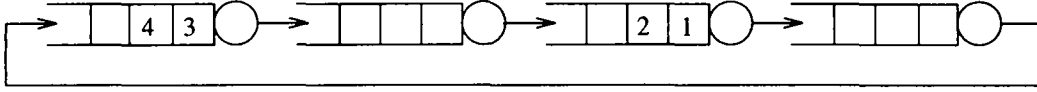


Figure 10: The model-1 network with one more station than that of Figure 1.

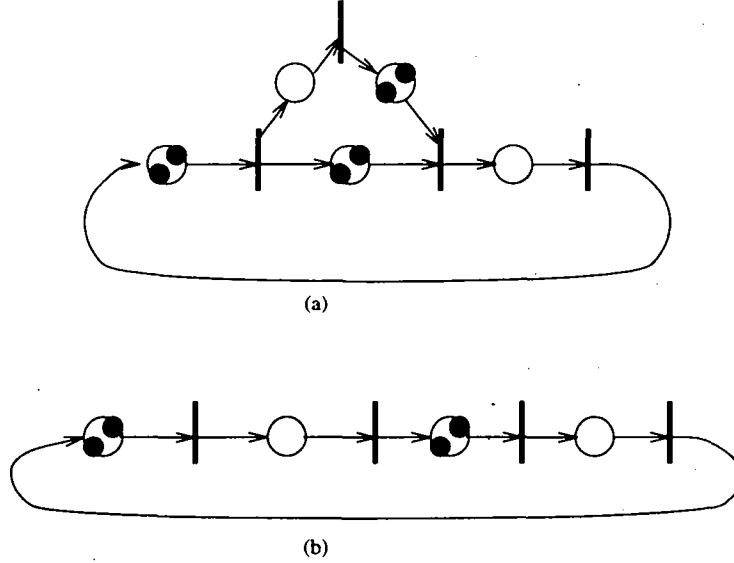


Figure 11: Stochastic marked graphs corresponding to the model-1 network of Figure 10. (a) Marked graph obtained by adding transitions and places to that of Figure 6. (b) Equivalent marked graph obtained by removing redundant places.

Let  $C_{n,k}$  (resp.  $\hat{C}_{n,k}$ ) be the  $k$ -th service completion time of server  $n$  in  $\mathcal{N}$  (resp.  $\hat{\mathcal{N}}$ ). Let  $l_n$  (resp.  $\hat{l}_n$ ) be the initial queue length of station  $n$  in  $\mathcal{N}$  (resp.  $\hat{\mathcal{N}}$ ). As we mentioned previously, the throughput of  $\mathcal{N}(N, M)$  and  $\hat{\mathcal{N}}(N+1, M+1)$  are independent of the initial queue lengths, we can assume without loss of generality that

$$\forall n, 1 \leq n \leq N : l_n = \hat{l}_n, \quad \text{and} \quad l_N = \hat{l}_N = \hat{l}_{N+1} = 1.$$

It is easily verified that the following recursive equations hold:

$$C_{n,k} = \max(C_{n,k-1}, C_{n-1,k-l_n}) + \sigma_{n,k}, \quad 1 \leq n \leq N, \quad k \geq 1, \quad (\text{B.3})$$

$$\hat{C}_{n,k} = \max(\hat{C}_{n,k-1}, \hat{C}_{n-1,k-\hat{l}_n}) + \hat{\sigma}_{n,k}, \quad 1 \leq n \leq N+1, \quad k \geq 1, \quad (\text{B.4})$$

where, by convention,  $C_{0,k} \equiv C_{N,k}$ ,  $\hat{C}_{0,k} \equiv \hat{C}_{N+1,k}$ , and if  $k \leq 0$ , then  $C_{n,k} = 0$  and  $\hat{C}_{n,k} = 0$ .

Let  $X_{n,k}$ ,  $1 \leq n \leq N+1$ ,  $k \geq 1$ , be the variables defined by the recursive equation

$$X_{n,k} = \max(X_{n,k-1}, X_{n-1,k-l_n}) + \sigma_{n,k}, \quad 1 \leq n \leq N+1, \quad k \geq 1, \quad (\text{B.5})$$

where, by convention,  $\sigma_{N+1,k} \equiv \sigma_{N,k}$ ,  $l_{N+1} \equiv l_N = 1$ ,  $X_{0,k} \equiv X_{N+1,k}$ , and if  $k \leq 0$ , then  $X_{n,k} = 0$ .

It is readily shown by induction that

$$X_{N+1,k} = X_{N,k}, \quad X_{n,k} = C_{n,k}, \quad 1 \leq n \leq N, \quad k \geq 1. \quad (\text{B.6})$$

Therefore, applying Jensen's inequality on conditional expectations in (B.4) implies

$$E[\hat{C}_{n,k}|U] \geq \max(E[\hat{C}_{n,k-1}|U], E[\hat{C}_{n-1,k-i_n}|U]) + E[\hat{\sigma}_{n,k}|U], \quad 1 \leq n \leq N+1, \quad k \geq 1, \quad (\text{B.7})$$

where  $U \stackrel{\text{def}}{=} \{\sigma_{n,k}^s, 1 \leq n \leq N\}_{k=1}^\infty$ . Replacing (B.1) and (B.2) in (B.7) yields

$$E[\hat{C}_{n,k}|U] \geq \max(E[\hat{C}_{n,k-1}|U], E[\hat{C}_{n-1,k-i_n}|U]) + \sigma_{n,k}, \quad 1 \leq n \leq N+1, \quad k \geq 1, \quad (\text{B.8})$$

By a simple induction proof using relations (B.5), (B.6) and (B.8), we can show that

$$E[\hat{C}_{n,k}|U] \geq C_{n,k} = E[C_{n,k}|U], \quad 1 \leq n \leq N, \quad k \geq 1. \quad (\text{B.9})$$

Unconditioning with respect to  $U$  implies that

$$E[\hat{C}_{n,k}] \geq E[C_{n,k}], \quad 1 \leq n \leq N, \quad k \geq 1. \quad (\text{B.10})$$

Hence,

$$\hat{\theta}(N+1, M+1) = \lim_{k \rightarrow \infty} \frac{k}{E[\hat{C}_{1,k}]} \leq \lim_{k \rightarrow \infty} \frac{k}{E[C_{1,k}]} = \theta(N, M).$$

The proof is thus completed. ■

Note that by Jensen's inequality, relation (B.9) implies the following increasing convex ordering

$$\{C_{n,k}, 1 \leq n \leq N\}_{k=1}^\infty \leq_{\text{icx}} \{\hat{C}_{n,k}, 1 \leq n \leq N\}_{k=1}^\infty. \quad (\text{B.11})$$

### Proof of Theorem 5.2.

We first show that insertion of stations and customers in a model-2 network results in additions of transitions and places in the corresponding marked graph. Indeed, insertion of a station results in an addition of a column of (two at each row) transitions and (two at each row) places. Insertion of a customer results in an addition of a row of (two at each column) transitions and (two at each column) places.

Such a fact is illustrated in Figure 13 which corresponds to the model-2 network of Figure 12. In the network of Figure 12, a new station with initially one customer is inserted on the right most of the stations of the network of Figure 2. In order to model this new station, we insert a column of “thick-bar” and “thin-bar” transitions as well as the places for connections in the marked graph of Figure 7 corresponding to the model-2 network of Figure 2. The two transitions and two places of each row in this new column represent the servers of the new station dedicated to existing customers. Recall that the “thick-bar” transitions represent the services, and the “thin-bar” transitions represent the departures and have zero firing time. In order to model the new customer, we insert a row of “thick-bar” and “thin-bar” transitions as well as the places for the connections. The two transitions and two places of each column in this new row represent the servers of each station dedicated to this new customer. Figure 13-(a) illustrates the resulting marked graph. Now, by removing the redundant places of this marked graph, we obtain the marked graph (Figure 13) corresponding to the network of Figure 12.

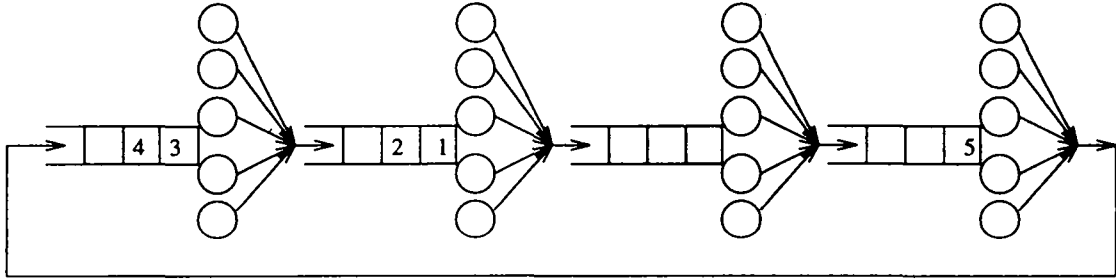


Figure 12: The model-2 network with one more station and one more customer than that of Figure 2.

The assertion of Theorem 5.2 becomes now a consequence of the properties of throughput monotonicity with respect to topology and with respect to timing established in [6] for stochastic marked graphs. ■



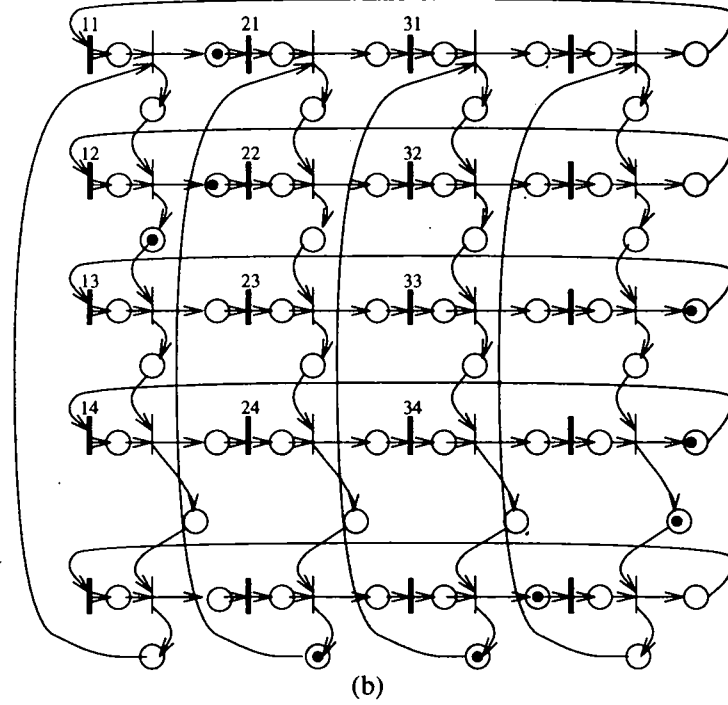
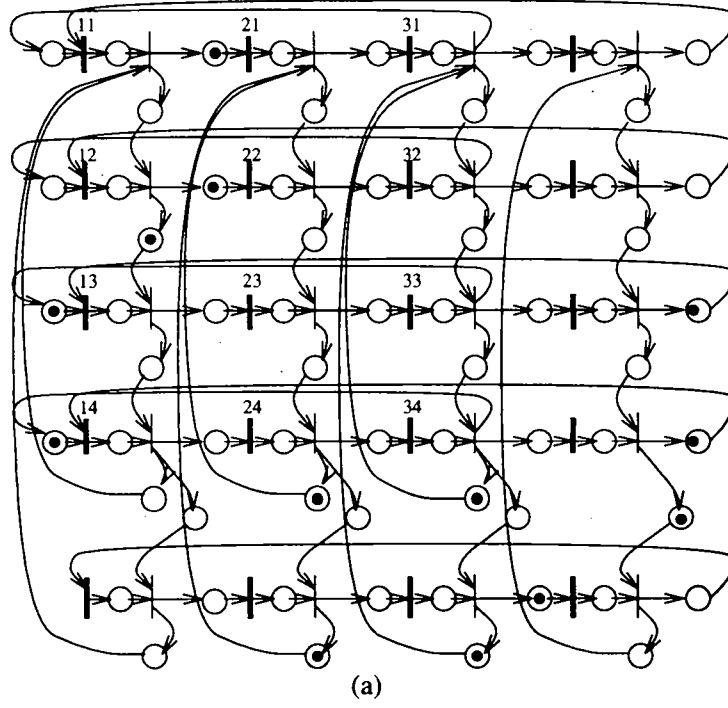


Figure 13: Stochastic marked graph corresponding to the model-2 network of Figure 12. (a) Marked graph obtained by adding transitions and places to that of Figure 7. (b) Equivalent marked graph obtained by removing redundant places in (a).

## References

- [1] I. F. Akyildiz, "Exact analysis of queueing networks with rejection blocking." In: *Queueing Networks with Blocking*, H.G. Perros and T. Altioek (Eds.), 1989, pp 19-29.
- [2] M. H. Ammar, S. B. Gershwin, "Equivalence relations in queueing models of fork/join queueing networks with blocking", *Performance Evaluation* **10** (1989), pp. 233-245.
- [3] F. Baccelli, P. Konstantopoulos, "Estimates of cycle times in stochastic Petri nets", in *Proc. of US-French Workshop on Applied Stochastic Analysis*, I. Karatzas, D. Ocone (Eds), Lecture Notes in Control and Information Sciences, **177**, 1991, pp. 1-20.
- [4] F. Baccelli, G. Cohen, G. J. Olsder, J.-P. Quadrat, *Synchronization and Linearity*, J. Wiley, 1992.
- [5] F. Baccelli, Z. Liu, "On a class of stochastic recursive sequences arising in queueing theory", *The Annals of Probability*, **20** (1992), pp. 350-374.
- [6] F. Baccelli, Z. Liu, "Comparison properties of stochastic decision free Petri nets", *IEEE Trans. on Automatic Control*, **37** (1992), pp. 1905-1920.
- [7] F. Baccelli, A. M. Makowski, "Queueing models for systems with synchronization constraints", *Proc. IEEE* **77** (1989), pp. 138-161.
- [8] O. J. Boxma, "Sojourn times in cyclic queues - the influence of the slowest server", In: *Computer Performance and Reliability*, G. Iazeolla, P.J. Courtois and O.J. Boxma (Eds.) North-Holland, 1988, pp 75-98.
- [9] J. P. Buzen, "Computational algorithms for closed queueing networks with exponential servers", *Comm. ACM* **16** (1973), pp. 527-531.
- [10] Y. Dallery, S. B. Gershwin, "Manufacturing flow line systems: A review of models and analytical results", Rapport de Recherche MASI, No. 91-18, 1991, Univ. of Paris 6.
- [11] Y. Dallery, Z. Liu, D. Towsley, "Equivalence, reversibility and symmetry properties in assembly/disassembly networks." Rapport de Recherche INRIA, No. 1267, 1990. To appear in the *Journal of the ACM*.
- [12] Y. Dallery, Z. Liu, D. Towsley, "Properties of fork/join queueing networks with blocking under various operating mechanisms", COINS Technical Report, TR 92-39, 1992.
- [13] E. Gelenbe, I. Mitran, *Analysis and Synthesis of Computer Systems*, Academic Press, 1980.
- [14] W. J. Gordon, G. F. Newell, "Cyclic queueing systems with restricted length queues", *Operations Research* **15** (1967), pp. 266-277.

- [15] S. Karlin, A. Novikoff, "Generalized convex inequalities", *Pacific J. Math.*, **13** (1963), pp. 1251-1279.
- [16] R. O. Onvural, "Survey of closed queueing networks with blocking", *ACM Computing Surveys* **22** (1990), pp. 83-122.
- [17] H. G. Perros, "A bibliography of papers on queueing networks with finite capacity queues", *Performance Evaluation* **10** (1989), pp. 255-260.
- [18] V. Rego, W. Szpankowski, "Closed-network duals of multiqueues with application to token-passing systems", *Computer Systems Science and Engineering*, **3** (1988), pp. 127-139.
- [19] R. Schassberger, H. Daduna, "The time for a round trip in a cycle of exponential queues", *J. ACM* **30** (1983), pp. 146-150.
- [20] V. Strassen, "The existence of probability measures with given marginals," *Ann. Math. Stat.*, Vol. 36, pp. 423-439, 1965.
- [21] R. Weber, "Scheduling and interchangeability in tandem queues," In *Scheduling Theory and Its Applications*, P. Chretienne et al. (Eds.), J. Wiley, 1993, to appear.



---

Unité de Recherche INRIA Sophia Antipolis  
2004, route des Lucioles - B.P. 93 - 06902 SOPHIA ANTIPOLIS Cedex (France)

Unité de Recherche INRIA Lorraine Technopôle de Nancy-Brabois - Campus Scientifique  
615, rue du Jardin Botanique - B.P. 101 - 54602 VILLERS LES NANCY Cedex (France)  
Unité de Recherche INRIA Rennes IRISA, Campus Universitaire de Beaulieu 35042 RENNES Cedex (France)  
Unité de Recherche INRIA Rhône-Alpes 46, avenue Félix Viallet - 38031 GRENOBLE Cedex (France)  
Unité de Recherche INRIA Rocquencourt Domaine de Voluceau - Rocquencourt - B.P. 105 - 78153 LE CHESNAY Cedex (France)

---

EDITEUR  
INRIA - Domaine de Voluceau - Rocquencourt - B.P. 105 - 78153 LE CHESNAY Cedex (France)

ISSN 0249 - 6399



★ R R - 2 1 1 5 ★